



**PRACTICAL INVESTIGATIONS ON BAYESIAN  
INVERSE PROBLEMS**

**By**

**Muzaffer Ege Alper**

*MSc. in Computer Engineering*

**Supervisor:**

**Alexandre Thiery**

**PHD IN STATISTICS AND APPLIED PROBABILITY**

**DEPARTMENT OF STATISTICS AND APPLIED PROBABILITY  
NATIONAL UNIVERSITY OF SINGAPORE**

**2016**

## Declaration

I hereby declare that this thesis is my original work and it has been written by me in its entirety. I have duly acknowledged all the sources of information which have been used in the thesis.

This thesis has also not been submitted for any degree in any university previously.



---

Muzaffer Ege Alper

November 2016

## Acknowledgments

It feels like I have been listening to Beckett's advice for quite some time. So why not now?

So long, thanks for all the fish and the kimchi fried rice.

## Abstract

Inverse problems make up a challenging and practically important class of inference problems. From an inferential perspective, the difficulty is in the ill-posed nature of the inverse problems, making least squares or maximum likelihood approaches inapplicable. Another difficulty is the computational complexity of the related forward problem, since usually one needs to solve one forward problem per likelihood evaluation. Tikhonov regularization is a class of methods that turns the problem into a well posed one while maintaining desirable statistical properties. Due to the fact that a Tikhonov regularization solution is obtained through maximization of a regularized fitness function, it is also computationally feasible and hence the standard method in the literature.

As the computational power increased, however, practitioners and researchers looked for better uncertainty quantification. The usual asymptotic confidence intervals gave way to full distributions using the Bayesian approach. This approach to inverse problems, while providing a full posterior distribution instead of a single point estimate as its answer, is also computationally much more expensive. Hence a great deal of research goes into finding computationally feasible methods of Bayesian inference. As is the case in the general Bayesian statistical literature, one may either opt to go for an “exact” solution using a Monte Carlo approach, or a fast approximate solution using Variational Bayes and other approximations. Aside from the complexity of likelihood evaluations, the Bayesian approach faces another problem: the difficulty of exploring very high dimensional distributions. One can try to speed up inference by having a more clever general purpose Monte Carlo method, or by trying to utilize problem specific features. On the other hand, recent research on “MCMC on functional spaces” provide the tools to work with very high dimensional distributions.

This thesis is made up of four main chapters. In the first two chapters, we review the literature on inverse problems and Monte Carlo methods. We put special emphasis on the functional space approach, which fits naturally into our programme of working in high dimensional inverse problems. We also attempt to describe alternative methods and, if possible, their relations to our proposed methods. We use the groundwater-flow dynamics as a basis to construct several inverse problems, which are later used as examples for our proposed methods. The third chapter describes our novel adaptive sequential Monte Carlo method and its application to the groundwater-flow problem. Here, we observe significant time-savings compared to previous SMC approaches. We also observe, however, that this method is still too slow to be used in practice. Therefore, in chapter four, we turn our attention to multi-resolution (also known as multi-level) methods. We begin by discussing a classical example in the literature, we clearly demonstrate the (asymptotically) reduced error per unit computation. We then describe our implementation of this idea and show the match of experimental results to the predictions of asymptotic theory. We end the thesis with our conclusions and observations on the practical properties of the discussed methods.

This thesis is mainly a practical one and hence does not contain any new theorems. However, it's important to know certain theorems to make sense of the resulting algorithms and get a sense of their behaviour. For that reason, we will state important ones when necessary but avoid giving the proofs. References to the proofs, of course, will be provided.

# Contents

<b>List of Figures</b>	<b>viii</b>
<b>List of Tables</b>	<b>ix</b>
<b>1 Inference in Inverse Problems</b>	<b>1</b>
1.1 Inverse Problems	1
1.1.1 Examples	3
1.2 Classical Inference in Inverse Problems	5
1.2.1 Basic Theory	5
1.3 Bayesian Inference in Inverse Problems	7
1.3.1 Prior modeling	8
1.3.2 Setting up the posterior	9
1.4 Groundwater Flow Problem	11
1.5 Numerical Solution of the Forward Problem: Finite Elements Method	14
<b>2 Monte Carlo Inference</b>	<b>17</b>
2.1 Markov Chain Monte Carlo Inference	17
2.1.1 Motivation and definition	17
2.1.2 Basic theory	18
2.1.3 MCMC in functional spaces	23
2.1.4 Tempered Monte Carlo methods	25
2.2 Sequential Monte Carlo Samplers	27
2.2.1 Motivation and definition	28
2.2.2 Curse of Dimensionality Discussions	31
2.2.3 A basic Central Limit Theorem for SMC	34
2.2.4 Feynman-Kac measures and some results	35
<b>3 Uni-Resolution Adaptive Sequential Monte Carlo Inference</b>	<b>38</b>

---

3.1	Problem Statement	38
3.2	Tempered Sequential Monte Carlo	39
3.3	Convergence Properties of Non-Adaptive SMC	40
3.4	Adaptive SMC	42
3.5	Numerical Results	44
3.5.1	Implementation details	44
3.5.2	Objects of inference	45
3.5.3	2D Results	46
3.5.4	3D Results	49
<b>4</b>	<b>Multi-Resolution Sequential Monte Carlo Inference</b>	<b>52</b>
4.1	Speeding up Monte Carlo Computations in Inverse Problems	52
4.1.1	Early Rejection	53
4.1.2	Delayed Acceptance	54
4.1.3	Surrogate Models	56
4.2	Multi-resolution Monte Carlo	59
4.3	Multi-resolution Monte Carlo Path Simulation	61
4.4	Multi-resolution Sequential Monte Carlo	64
4.5	Results	66
<b>5</b>	<b>Conclusion</b>	<b>69</b>
	<b>References</b>	<b>71</b>

# List of Figures

1.1	1D FEM segments and linear basis	16
2.1	Acceptance rates for the symmetric random walk proposal, when the target is the same as the prior	24
2.2	$d = 10$	33
2.3	$d = 40$	33
2.4	$d = 100$	33
2.5	$d = 400$	33
2.6	$\max(w)$ plots for different dimensions $d$ .	33
3.1	Standard SMC Samplers. $M_{\text{thres}} \in \{1, \dots, M\}$ is a user defined parameter.	40
3.2	Histograms of a single component of SMC particles at different iterations of the algorithm corresponding to the adaptively selected temperatures, together with the target tempered posteriors	43
3.3	Six permeability field samples drawn from the prior	47
3.4	Numerical consistency checks for the sequence of experiments with 4,16,36,64 and 100 observations	48
3.5	True Permeability Field	48
3.6	Estimated Permeability Field	48
3.7	An estimated permeability field and the corresponding true field	48
3.8	Posterior marginal density estimates for two low and one high frequency coefficients in the 2D case	49
3.9	SMC Performance for 3D Example.	50
3.10	Posterior marginal density estimates for two low and one high frequency coefficients in the 3D case	51
4.1	log-Cost vs. log-Error plot for three methods	67



## List of Tables

- 3.1 Parameter values used for the 2D experiments. Between 5 and 1000 steps are allowed for the iterates of the MCMC kernels. The frequency cutoff determines the level of discretization of the permeability field. Finite elements d.o.f. denotes the number of finite elements used in the numerical solution of the elliptic PDE, higher values indicate better approximation at the expense of computational resources. 46
- 3.2 Parameter values used for the 3D experiment. Between 5 and 200 steps are allowed for the iterates of the MCMC kernels. 49

# CHAPTER 1

## Inference in Inverse Problems

In this chapter, we present the basics of inverse problems and Bayesian inference that will be used repeatedly throughout this thesis. There is an abundance of models for linear and nonlinear inverse problems which address problems arising in diverse practical fields. A quick overview of the literature produces as diverse applications as; thermal analysis of re-entry space shuttles [BBC85], electrocardiography [Gul05], corrosion detection and modeling [Ing97], computerized tomography [Nat01], positron emission tomography [SV82] among many others [Ron08, AS09, WA11]. A central challenge in the inverse problems is to overcome the ill-posedness of the inverse map, where a small change in the observables make a big impact on the solution. Our goal in this chapter is to hint at the diversity of inverse problems through examples, illustrate the challenges of their ill-posed nature and describe the common approaches to overcome this issue. In particular, we will briefly describe Tikhonov regularization based approaches and contrast this to our Bayesian approach. The discussions of this chapter will act as a basis for the rest of this thesis.

We reiterate that none of the theorems stated here are new and hence they are stated without proofs.

### 1.1 Inverse Problems

In this section, we will discuss the essentials of the class of problems called “Inverse problems”. The statistical and computational features of these problems set them apart from other inference problems. We will briefly describe such features and provide simple examples when possible. This section is heavily influenced by the work of Ito et.al. [IJ14].

The objective of the inverse problems is to construct a stable inverse map from observables to unknowns, one with hopefully good statistical properties. The observables are assumed to be generated by a forward map, mapping the parameter space to the observable space, contaminated by observation noise. A major challenge in inverse problems is the ill-posedness of the inverse map. A well posed problem, as defined by French mathematician Jacques Hadamard satisfies these properties;

- a) **(Existence)** There exists at least one solution.
- b) **(Uniqueness)** There is at most one solution.
- c) **(Stability)** The solution depends continuously on data.

Historically ill-posed problems were considered uninteresting and irrelevant. However, with such inverse problems appearing in broad areas like medical imaging, weather prediction, petroleum source detection and so on, it is impossible to ignore them. An early successful solution to such problems is Tikhonov Regularization, where you maximize a regularized fitness function. This approach manages to alleviate the aforementioned issues like possible non-existence, non-uniqueness or instability.

In applications, the observables, usually few in number, only provide partial information of the unknowns which are usually very high dimensional. Hence, point estimates can be very misleading, making uncertainty quantification crucial. A modern approach to quantify uncertainty employs the Bayesian stance, outputting the full posterior distribution, instead of just a point estimate. Prior information is often naturally available, making this approach even more fitting. We leave the details of the Bayesian approach to the next section.

We now give examples for ill-posedness of the inverse map. The non-existence and non-uniqueness cases turn out to be the easiest to understand, but also the least serious problems. One can obtain a non-existence case by considering observation noise that can possibly leave the observables outside the range of the forward map. Likewise, solution to an underdetermined linear system will be non-unique.

Consider the following simple problem as an example to an unstable inverse map. Suppose that a particle is moving in a straight line, and that the motion is caused by the force  $q(t)$ ,

depending on time. Assume the particle is at origin with no motion at  $t = 0$ . The motion of the particle is described by  $u(t)$  satisfying the following differential equation:

$$\begin{aligned} \frac{d^2 u}{dt^2} &= q(t); & \text{for } t \in [0, T] \\ u(0) &= 0, & \frac{du}{dt}(0) = 0, \end{aligned}$$

where  $u(t)$  is the position of the particle at time  $t$ . The goal of the inverse problem, then, is to construct the force  $q(t)$  from limited (or complete) information of  $u(t)$ . To understand the stability of this map, we investigate how the solution  $q(t)$  behaves when we make a small perturbation to  $u(t)$ . If we take  $u_n(t) = u(t) + \frac{1}{n} \cos(nt)$ , the corresponding solutions become  $q_n(t) = q(t) - n \cos(nt)$ . We now observe that  $\|u - u_n\|_\infty \rightarrow 0$  as  $n \rightarrow \infty$  while,  $\|q - q_n\|_\infty \rightarrow \infty$ , i.e. very small perturbations in observations can cause arbitrarily large differences in the solutions.

### 1.1.1 Examples

To make things more concrete, we use the general elliptic problem as a template. The forward problem looks like this:

$$\nabla \cdot (a(x) \nabla u) + \mathbf{b}(x) \cdot \nabla u + p(x)u = f(x); \quad \mathbf{x} \text{ in } \Omega \quad (1.1)$$

$u$  above is the solution of the forward map, i.e. the solution of the Partial Differential Equation (PDE). The parameters are as follows;  $a$  is conductivity/diffusivity,  $b$  is the convection coefficient,  $p$  is potential and  $f$  is the source/sink term. In practice, one usually picks one of the parameters and assumes the rest fixed (in particular, they may be 0). The goal, then, is to make inference on this chosen parameter given some observations on  $u$ . In this way, one can come up with many inverse problems, like conductivity-inversion, source-inversion and so on. The template elliptic problem given above does not involve any time parameter and hence usually corresponds to steady-state behaviour of some system. Let us consider two typical inverse problems generated from this template. We will see a third example, which is the focus of this thesis, later in this chapter.

### Cauchy problem

The Cauchy problem is formulated as follows. Let  $\Gamma_c$  and  $\Gamma_i = \Gamma \setminus \Gamma_c$  be two disjoint parts of the boundary  $\Gamma$ , referring to the observed and non-observed parts. Then, given Cauchy observations  $(g, h)$  on the boundary  $\Gamma_c$ , the goal is to infer  $u$  on the boundary  $\Gamma_i$ , i.e.,

$$-\nabla \cdot (a(x)\nabla u) = 0; \quad \text{in } \Omega$$

$$u = g; \quad \text{on } \Gamma_c$$

$$a \frac{\partial u}{\partial n} = h; \quad \text{on } \Gamma_c$$

where  $\frac{\partial}{\partial n}$  denotes partial derivative with respect to the boundary normal. One application of this inverse problem is the thermal analysis of re-entry space shuttles. Here, one can measure the temperature and heat flux on the inner surface of the shuttle, and one is interested in the flux on the outer surface, which is not directly accessible [IJ14].

### Inverse source problem

A second classical problem can be generated from Eq. 1.1 by inverting the source term  $f$ , i.e.,

$$-\Delta u = f$$

$$u = g; \text{ and } \frac{\partial u}{\partial n} = h \quad \text{on } \Gamma.$$

An exemplary application is electroencephalography, where brain's spontaneous electrical activities are recorded from electrodes placed on the scalp. Retrieving source term from the Cauchy data is not unique and one usually requires additional sources of information (such as a prior, or certain constraints) to be able to make inference.

## 1.2 Classical Inference in Inverse Problems

In this section, we will briefly describe the Tikhonov regularization method and the basics of the associated theory, including well-posedness and the consistency of the solution. As before, we consider possibly nonlinear operator equations:

$$G(u) = g^{t1}$$

where the operator  $G : \mathcal{X} \rightarrow \mathcal{Y}$  ( $\mathcal{X}, \mathcal{Y}$  are Hilbert spaces) is called the “forward operator”, possibly a composition of the solution operator to an underlying PDE and a selection/sampling operator for pointwise observations.  $g^t$  denotes the true observations, i.e. without noise contamination. As before, in practice we only have access to noisy data  $g^\delta$ , whose accuracy with respect to the true data is measured by the noise level  $\delta$ :

$$\delta = \|g^t - g^\delta\|.$$

The classical approach for obtaining a well-behaved, approximate solution is Tikhonov regularization, which is the minimizer of Tikhonov functional:

$$J_\alpha(u) = \|G(u) - g^\delta\|_p^p + \alpha\psi(u),$$

where the first term incorporates the information from observed data, while the second term serves to regularize the solution and to incorporate prior information.

### 1.2.1 Basic Theory

We will now briefly discuss the well-posedness, that is the existence, stability and consistency of Tikhonov solutions. The following set of assumptions are central [IJ14]:

**Assumption 1.** *The operator  $G : \mathcal{X} \rightarrow \mathcal{Y}$ , the non-negative regularization functional  $\psi$  and the corresponding Tikhonov functional  $J_\alpha(u)$  satisfy;*

---

<sup>1</sup>note that  $u$  denotes the parameter now, as opposed to the solution of a PDE as in the previous sections

1.  $J_\alpha(u_n) \rightarrow \infty$  as  $\|u_n\| \rightarrow \infty$ .
2.  $u_n \rightarrow u^*$  weakly in  $X$  implies  $G(u_n) \rightarrow G(u^*)$  weakly in  $Y$ .
3. The functional  $\psi$  is proper convex and weakly lower semicontinuous.

We now state the existence result as in [IJ14].

**Theorem 1.** *Let Assumption 1 hold. Then for every  $\alpha > 0$ , there exists a minimizer to  $J_\alpha$  <sup>2</sup>.*

Define  $u_\alpha^\delta$  as the minimizer to  $J_\alpha(u)$  with noise level  $\delta$ . As discussed before, the critical issue with inverse problems is the stability of the solution. Therefore, we now turn to the stable dependence of  $u_\alpha^\delta$  to perturbations in the data  $g^\delta$ .

**Theorem 2.** *Let Assumption 1 hold. Let  $\{g_n\}$  be a sequence converging to  $g^\delta$  (true data) in  $Y$ , and  $\{u_n\}$  be a sequence of minimizers to  $J_\alpha$  with  $g_n$  in place of  $g^\delta$ . Then the sequence  $\{u_n\}$  contains a weakly convergent subsequence, and the limit is a minimizer to the functional  $J_\alpha$ .*

Finally, a natural requirement for a successful inference method is to have convergence to the true solution as the noise level diminishes. In the case of inverse problems however, the problem need not have a single solution even under perfect (no-noise) conditions. This is in particular true when the forward operator  $G$  is the composition of sampling and PDE solution operators, which is the case of interest in this thesis. Therefore, we need a rigorous notion of “true solution”. For this purpose, we define the “ $\psi$ -minimizing solution” as the element(s)  $u^t$  that satisfies:

$$\psi(u^t) \leq \psi(u); \quad \forall u \in \{u : G(u) = g^t\}.$$

The following theorem [IJ14] states the convergence of Tikhonov solutions to the  $\psi$ -minimizing solution when the error level diminishes, when the regularization parameter  $\alpha$  is appropriately selected.

**Theorem 3.** *Let Assumption 1 hold. Let the sequence  $\{\delta_n\}$  be convergent to 0, and  $g^{\delta_n}$*

---

<sup>2</sup>not necessarily unique!

satisfy  $\|g^t - g^{\delta_n}\| = \delta_n$ . Choose the parameter  $\alpha(\delta)$  s.t.

$$\lim_{\delta \rightarrow 0} \alpha(\delta) = 0 \quad \text{and} \quad \lim_{\delta \rightarrow 0} \frac{\delta^p}{\alpha(\delta)} = 0.$$

Let  $\{u_{\alpha(\delta)}^{\delta_n}\}$  be a sequence of minimizers to  $J_{\alpha(\delta)}$  with  $g^{\delta_n}$  in place of  $g^t$ . Then it contains a subsequence converging weakly to  $\psi$ -minimizing solution.

These conclude the basic properties of the Tikhonov solution. There is, of course, a rich literature on this subject with deeper results. We refer to, e.g., [IJ14, BPR07, Neu98, Zhd93, Ten01].

### 1.3 Bayesian Inference in Inverse Problems

In this section, we shall discuss the essential components of Bayesian inference: setting up a prior and using the Bayes theorem. We will consider the parameter, to be denoted as  $u$  whether it is permeability or forcing, to lie in a general separable Hilbert space. Such generality will allow us to consider function-valued parameters (e.g. fields) as well as the usual vector-valued case in the same framework. This section makes heavy use of the discussion provided in the paper by Dashti et.al. [DS13].

Before we begin, let us briefly motivate the advantage of considering the function-space setting. As argued in the important paper of Beskos & Stuart [BS09], one can consider the sampling performance of a Monte Carlo method in the high dimension either by taking the limit of the acceptance rates w.r.t. increasing dimensionality or by considering the Monte Carlo method directly in the functional (“infinite dimensional”) space. The latter has the advantage that the practical high-but-finite dimensional problem is just a subset of the infinite dimensional problem, therefore any property of the latter will carryover to the former. In particular, if the Monte Carlo method is well defined in the function-space setting, e.g. having positive acceptance rate for MCMC, than it will behave well for any finite dimensional case too. The authors of the aforementioned paper use this idea to build a MCMC kernel that is well defined in the function space setting which will therefore behave



well no matter how large you set the dimensionality.

### 1.3.1 Prior modeling

We assume our parameter lies in a separable Hilbert space  $\mathcal{X}$ . We use countable infinite sequences to model the parameter function:

$$u = \phi_0 + \sum_{j=1}^{\infty} c_j \phi_j.$$

By randomizing  $c_j$ , we can create a real valued random function. We consider the domain  $\mathcal{D}$  of  $\phi_j$ s as either an open bounded subset of  $\mathbb{R}^d$  or as the torus  $\mathcal{T}^d$ . We set  $c_j = \gamma_j u_j$ , where  $\gamma_j$  is the deterministic part and  $u_j$  is the random part. We can then choose  $\gamma_j$  s.t. we get a convergent series. In our work, we use Fourier basis for  $\phi_j$ s and we pick  $\gamma_j = \frac{1}{\|k(j)\|_{\infty}^{d+0.001}}$ , where  $k(j)$  maps to the frequency of the  $j$ th coefficient (i.e.  $k(j)$  is a vector, and hence  $\|k(j)\|_{\infty}$  is its largest component).

#### Uniform case

A uniform prior can be set-up by choosing  $\gamma = (\gamma_j)_j \in l^1$  and  $u_j \sim U[-1, 1]$ . Here  $l^1$  is the space of absolutely summable real sequences and  $U[a, b]$  is the uniform distribution in range  $U[a, b]$ . Now assume there exists constants  $\phi_{min}, \phi_{max}, \delta$  s.t.

$$\text{ess inf}_x \phi(x) \geq \phi_{min} > 0$$

$$\text{ess sup}_x \phi(x) \leq \phi_{max}$$

$$\|\gamma\|_{l^1} = \frac{\delta}{1 + \delta} \phi_{min},$$

where  $\|x\|_{l^1} = \sum_i |x_i|$ . We obtain the following property;

$$\frac{1}{1 + \delta} \phi_{min} \leq u(x) \leq \phi_{max} + \frac{\delta}{1 + \delta} \phi_{min} \quad a.e..$$

This setup may not look very realistic or applicable, but it turns out that it is widely employed. One example is the 2012 paper by Hoang et.al. [HSS12]. A particular advantage of this scheme is that it lends itself more easily to analysis, hence it is employed more often in analysis papers. The alternative to be described next is more often used in applications.

### Log-Gaussian case

We now consider  $u$  as the logarithm of the parameter of interest. This approach has the advantage of guaranteeing that the original parameter is positive and hence we do not need the strict assumptions made in the previous section. Assume  $\{\phi_j\}_j$  is an orthogonal basis for the separable Hilbert space  $\mathcal{X}$ . We take  $u_j \sim N(0, 1)$ , thus having  $u \sim N(\phi_0, C)$ , where the covariance operator  $C$  depends on the choice of  $\gamma$ .  $\phi_0$  is usually taken as 0 in applications. We have the following theorem [DS13].

**Theorem 4.** *Assume  $\gamma_j = O(j^{-\frac{s}{d}})$ . Then the sequence of partial sums  $u^N = \sum_{i=1}^N c_i \phi_i$  is Cauchy in  $L^2(H^t)$  for  $t < s - \frac{s}{d}$ .*

We remind that  $L^2(H^t)$  has the norm  $\|u\|_{H^t} = \sum j^{\frac{2t}{d}} |u_j|^2$  and  $H^t$  is the (Hilbert) space of fields with finite  $\|\cdot\|_{H^t}$  norms.

Before we close this subsection on priors, we should point out that there are deeper results on the regularity of the generated functions than those stated here [DS13]. However, for our purposes, these suffice.

### 1.3.2 Setting up the posterior

The goal of this subsection is to give the conditions under which we have a Bayes theorem for the function-space setting so that the posterior is absolutely continuous w.r.t. the prior. The absolute continuity is not crucial (see: [HD12]), but it's a natural requirement since it implies that any almost sure property of the prior will carry-over to the posterior.

Let  $\mathcal{X}$  and  $\mathcal{Y}$  denote a pair of measurable spaces (with appropriate sigma-algebras) and let

$\pi, \nu$  be probability measures on  $\mathcal{X} \times \mathcal{Y}$ . Assume  $\nu \ll \pi$ , i.e. there exists  $\phi \in L^1_\pi$  s.t.

$$\frac{d\nu}{d\pi}(x, y) = \phi(x, y).$$

**Theorem 5.** [DS13] Assume  $\pi^y$  exists s.t.  $\pi^y(dx)\pi(dy) = \pi(dx, dy)$ . Then there also exists  $\nu^y$  s.t.  $\nu^y(dx)\nu(dy) = \nu(dx, dy)$ <sup>3</sup>. Furthermore, assume  $c(y) := \int_X \phi(x, y)d\pi^y(x) > 0$ , then  $\nu^y \ll \pi^y$  and,

$$\frac{d\nu^y}{d\pi^y}(x) = \frac{1}{c(y)}\phi(x, y).$$

Now assume the spaces  $\mathcal{X}, \mathcal{Y}$  are also separable Banach spaces and let  $G : \mathcal{X} \rightarrow \mathcal{Y}$  a measurable mapping. We want to formulate a Bayes theorem for problems of the sort:

$$y = G(u) + \epsilon$$

where  $\epsilon$  denotes noise and  $u \in \mathcal{X}$  is the parameter of interest.

Let  $(u, y)$  be a random variable, specified using;

a) prior:  $u \sim \mu_0$  (a measure on  $\mathcal{X}$ )

b) noise:  $\epsilon \sim Q_0$  (measure on  $\mathcal{Y}$ ) and  $\epsilon \perp u$

Hence the random variable  $y|u$  is distributed according to  $Q_u$ : the translate of  $Q_0$  by  $G(u)$ .

Assume  $Q_u \ll Q_0$  so that,

$$\frac{dQ_u}{dQ_0}(y) = \exp(-\Phi(u, y))$$

where, for a given instance of data  $y$ ,  $-\Phi(u, y)$  is termed the log-likelihood. Now define  $\nu_0(du, dy) := \mu_0(du)Q_0(dy)$  and assume  $\Phi(u, y)$  is  $\nu_0$  measurable. Then  $(u, y)$  is distributed according to  $\nu(du, dy) = \mu_0(du)Q_u(dy)$  and,

$$\frac{d\nu}{d\nu_0}(u, y) = \exp(-\Phi(u, y)).$$

**Theorem 6.** [DS13] Assume  $Z := \int_X \exp(-\Phi(u, y))\mu_0(du) > 0$ . Then  $\nu^y(du)$  exists and

---

<sup>3</sup> $\pi(dy) = \int_x \pi(dx, dy)$  and  $\nu(dy) = \int_x \nu(dx, dy)$

$\nu^y \ll \mu_0$ , additionally;

$$\frac{d\nu^y}{d\mu_0}(u) = \frac{1}{Z} \exp(-\Phi(u, y)).$$

The proof of the above theorem is essentially an application of Theorem 5 with  $\pi$  replaced with  $\nu_0$ ,  $\phi(x, y) = \exp(-\Phi(x, y))$ ,  $(x, y) \rightarrow (u, y)$  and  $\pi^y = \mu_0$ .

From the above discussion, we gather that there are three essential steps to apply this general Bayes theorem:

- 1) Ensure  $Q_u \ll Q_0$
- 2) Ensure  $\Phi$  is  $\nu_0$  measurable.
- 3) Ensure  $Z := \int_X \exp(-\Phi(u, y)) \mu_0(du) > 0$

We will show an application of this procedure under the uniform prior setup for the groundwater flow problem defined in the next section.

## 1.4 Groundwater Flow Problem

In this section, we will define the groundwater flow problem and describe its important properties such as the continuity of the forward map w.r.t. the parameter. We will end the discussion by showing how we can apply the Bayes theorem of the previous section in this setting, using uniform prior. The problem can be stated as follows; let  $p(\kappa)$  (“pressure field”) behave according to:

$$-\nabla \cdot (\kappa \nabla p) = f; \quad x \in D$$

$$p = 0; \quad x \in \partial D$$

where  $D$  is an open, bounded subset of  $\mathbb{R}^d$  and  $\partial D$  is its boundary. The inverse problem is to find  $\kappa$  given some set of linear bounded functionals of  $p$ , assuming  $f$  is known. A different inverse problem can be set-up by considering forcing ( $f$ ) inversion, with  $\kappa$  known. This second problem is much easier to deal with as it’s linear, and we shall consider it as a toy problem for the multi-resolution Monte Carlo approach. In the rest of this section, we will discuss only the permeability ( $\kappa$ ) inversion.

The corresponding weak formulation reads;

$$-\int_D \kappa \nabla p \cdot \nabla v dx = \int_D f v dx$$

where the so-called test function  $v \in V = (H_0^1(D), \langle \cdot, \cdot \rangle, \|\cdot\|)$ ,  $p \in H^1(L^2(D))$ . Recall that  $H_0^1(D)$  and  $H^1(L^2(D))$  are instances of what is called a Sobolev space. In general, a Sobolev space  $H^i(L^p)$  is a Banach space where the norm is the sum of  $L^p$  norms of the function and its derivatives up to  $i$ .  $H_0^1 \subset H^1$  is the space of functions in  $H^1$  that vanish at the boundary. We make the following common assumptions;

A1.  $\text{ess inf}_x \kappa(x) = \kappa_{\min} > 0$

A2.  $f \in V^*$ , i.e. the dual space of  $V$

Under these conditions, we have the following theorem;

**Theorem 7.** [DS13] *There is a unique weak solution to the above forward problem that satisfies;*

$$\|p\|_V \leq \|f\|_{V^*} / \kappa_{\min}.$$

This result is in turn used to prove the following continuity theorem;

**Theorem 8.** [DS13] *For  $i = 1, 2$ , let;*

$$-\nabla \cdot (\kappa_i \nabla p_i) = f, \quad x \in D$$

$$p_i = 0, \quad x \in \partial D.$$

Then,

$$\|p_1 - p_2\|_V \leq \frac{1}{\kappa_{\min}^2} \|f\|_{V^*} \|\kappa_1 - \kappa_2\|_{L^\infty}$$

with,

$$\kappa_{\min} := \text{ess inf}_x \kappa_1 \wedge \text{ess inf}_x \kappa_2 > 0.$$

Such continuity theorems are used in proofs of estimator stability [IJ14], posterior consistency [Vol13] and error analysis [HSS12].

We now move on to the application of Bayes theorem in this problem, assuming the uniform prior explained in the previous section. We begin by checking that the prior correctly assigns measure 1 to the admissible parameter space. We know from our previous discussion that;

$$\frac{1}{1+\delta}\phi_{min} \leq \kappa \leq \phi_{max} + \frac{\delta}{1+\delta}\phi_{min} \quad a.e.$$

and hence the admissible parameter space  $\{\kappa \in X | \text{essinf}_x \kappa(x) > 0\}$  has measure 1.

We then have to find  $\Phi(.,.)$  and check that it's appropriately measurable. The observations  $\{l_j\}_j$  are linear bounded functionals on  $V$ , hence  $l_j \in V^*$ . Define  $G_j(\kappa) = l_j(p(\kappa))$  and  $G(\kappa) := (G_1, G_2, \dots, G_J)$ , with  $y = G(\kappa^t)$  representing the vector of  $J$  observations under the true permeability  $\kappa^t$ . It is important to note that observations themselves are real-vectors, i.e. have finite dimensions. The function-valued observations can be of theoretical interest, but is of no concern to us for our practical investigations.  $\Phi$  is defined by the following derivative

$$\frac{dQ_u}{dQ_0}(\kappa) = \exp(-\Phi(\kappa, y)).$$

Assuming  $\epsilon \sim N(0, I)$ , we get

$$\Phi(\kappa, y) = 0.5|y - G(\kappa)|^2 - 0.5|y|^2.$$

Using Theorem 8, we conclude that  $\Phi$  is well defined. Measurability follows immediately from basic properties of measurable functions.

Finally we have to show that the normalizing constant is positive, i.e.  $Z := \int_X \exp(-\Phi(\kappa, y))\mu_0(d\kappa) > 0$ . We have shown that  $\kappa$  is bounded, therefore  $G$  is bounded in  $\mathfrak{R}^J$  and  $y$  is finite a.s. We arrive at the conclusion that  $\Phi < M(y) < \infty$  a.s. . Therefore,  $Z := \int_X \exp(-\Phi(\kappa, y))\mu_0(d\kappa) \geq \int_X \exp(-M)\mu_0(d\kappa) = \exp(-M) > 0$ .

We note that the Bayes theorem for the Gaussian prior case is done using the same steps, but is more convoluted. We refer the reader to Stuart et.al. [DS13].

## 1.5 Numerical Solution of the Forward Problem: Finite Elements Method

Analytical solution to the varying coefficient nonhomogeneous Laplace problem (e.g. the groundwater flow problem of the previous section) does not exist. Among the available methods, we employ Finite Elements Method (FEM) with polynomial basis functions. This is also known as the Galerkin approximation. We begin by expressing the problem in its weak formulation as before;

$$-\int_D \kappa \nabla p \cdot \nabla v dx = \int_D f v dx$$

where  $v$ , the test function, is defined as in the previous section. A general way of expressing such equations is by using bilinear form  $A(.,.)$  and linear form  $F(.)$  :

$$A(p, v) = F(v)$$

where  $A(p, v) = \int_D \kappa \nabla p \cdot \nabla v dx$  and  $F(v) = \int_D f v dx$ . A bilinear form  $A$  is called coercive if there exists an  $\alpha$ , s.t. for every  $u, v \in V$ ;

$$A(p, v) \geq \alpha \int_D \nabla p \cdot \nabla v dx.$$

In our setup, this is trivially true with  $\alpha = \kappa_{min}$ , where  $\text{ess inf}_x \kappa(x) = \kappa_{min} > 0$  (Assumption A1 of the previous section). This property, together with continuity of  $A$  and  $F$  imply a unique solution to the forward problem. We will see that these also imply the convergence of Galerkin approximations, to be defined next.

Assume that we have a family  $V_N$  of finite dimensional subspaces of  $V$ . Then the Galerkin approximation  $p_N \in V_N$  is defined by

$$A(p_N, v) = F(v)$$

for all  $v \in V_N$ . By the above discussion, the coercivity of  $A$  also holds true for  $V_n$ , hence the

Galerkin approximation exists uniquely.

Let us now consider the relationship of Galerkin approximation to the best approximation in  $V_N$ . The error of the best approximation is  $\inf_{q \in V_N} \|q - p\|_V$ , which is clearly less than or equal to  $\|p_N - p\|_V$ . This implies that if we have a convergent Galerkin approximation, we also have  $\inf_{q \in V_N} \|q - p\|_V \rightarrow 0$  as  $N \rightarrow \infty$ . Cea's lemma establishes the reverse statement:

**Lemma 1.** *If  $A$  is continuous and coercive then,*

$$\|p_N - p\|_V \leq C \inf_{q \in V_N} \|q - p\|_V.$$

The lemma says that the Galerkin solution is like the best approximation in  $V_N$  up to a constant that does not depend on  $N$ . For this reason, Galerkin approximation is sometimes called Galerkin projection. We also have the following result;

**Theorem 9.** *[Dur] If  $A$  is continuous and coercive and the spaces  $V_N$  are such that  $\inf_{q \in V_N} \|q - p\|_V \rightarrow 0$  as  $N \rightarrow \infty$ , then  $\lim_{N \rightarrow \infty} u_N = u$ .*

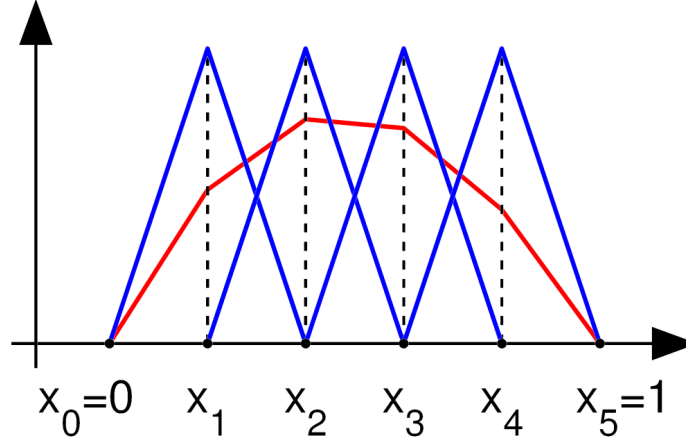
The next question is to construct “good” approximation subspaces  $V_N$  of  $V$ , the space where the exact solutions belongs. The Finite Element Method (FEM) provides a systematic way of constructing such spaces. The domain is divided into finite number of elements and polynomial basis functions inside these elements are used to construct the approximate  $p$ . We will explain the method in the 1D domain, i.e.  $D \in \mathfrak{R}$ .

Assume that we have a segmentation  $\mathcal{S} = \{S\}$  of  $D \in \mathfrak{R}$ , i.e.  $D = \bigcup_{S \in \mathcal{S}} S$ . Intersection of line segments  $S$  should contain only a common point, i.e. the segments are not allowed to overlap. Given a natural number  $k$ , we associate with  $\mathcal{S}$  the space  $V^k(\mathcal{S})$  of continuous piecewise polynomials of order  $k$ . We observe that  $V^k(\mathcal{S}) \subset V$ . Recall that we are working with homogeneous Dirichlet boundary conditions. A natural way to introduce this condition to FEM is to work with the subset  $V_0^k(\mathcal{S}) \subset V^k(\mathcal{S})$  of functions that vanish at the boundary. Therefore, we can define the FEM approximation  $p_{\mathcal{S}}$  as the exact solution of

$$A(p_{\mathcal{S}}, v) = L(v); \forall v \in V_0^k(\mathcal{S}).$$



Figure 1.1 illustrates a segmentation and the corresponding linear basis.



**Figure 1.1:** 1D FEM segments and linear basis

In order to apply Theorem 9, we need to show that these approximating spaces are “good enough”, i.e.  $\inf_{q \in V_N} \|q - p\|_V \rightarrow 0$  as  $N \rightarrow \infty$ . Equivalently, it suffices to show that there exist a good sequence of approximations in  $V_N$ , i.e. some  $q_N^*$  s.t.  $\|q_N^* - p\|_V \rightarrow 0$  as  $N \rightarrow \infty$ . Such convergence results are readily available in the polynomial approximation literature, we refer to Chapter 4 of [SS02].

Consider the piecewise linear basis functions  $\phi_k$ , centered around the intersection points of segments, called nodes. We choose  $\phi_k(x_k) = 1$  and  $\phi_k(x_j) = 0$  for  $j \neq k$ , where  $x_i$  are the nodal points. These basis functions form our subspace  $V_k$ . Let us express  $p$ ,  $v$  and  $f$  in these bases,  $p = \sum_i p_i \phi_i$ ,  $v = \sum_i v_i \phi_i$  and  $f = \sum_i f_i \phi_i$ . We can now re-express our approximate problem using these,

$$-L\mathbf{p} = M\mathbf{f}$$

where  $L_{ij} = \int \frac{d}{dx} \phi_i \frac{d}{dx} \phi_j dx$  and  $M_{ij} = \int \phi_i \phi_j dx$ . Observe that these matrices are analytically available. In addition, their size depends on the number of nodes, which in turn depends on the number of segments. Finally, a crucial feature of these matrices is that they are very sparse, since the bases are local. This means that it is possible to use efficient sparse solvers to obtain the solution.

# CHAPTER 2

## Monte Carlo Inference

An important requirement in applied inverse problems is uncertainty quantification (UQ). Bayesian paradigm provides a natural way to obtain UQ through posterior distributions. However, computing these distributions and extracting required information from them are non-trivial tasks. Bayesian inverse problems literature provides two answers: deterministic approximations [[FWA<sup>+</sup>11](#), [BBG<sup>+</sup>11](#), [BTBG<sup>+</sup>12](#)] and Monte Carlo inference. In this chapter, we will discuss the basics of the latter.

### 2.1 Markov Chain Monte Carlo Inference

Among the Monte Carlo family of inference tools, Markov Chain Monte Carlo (MCMC) methods are arguably the most popular. Likewise, a greater part of the Bayesian inverse problems literature focus on MCMC based inference. In this section, we will briefly motivate and define this class of methods, as well as giving a short overview of the associated theory. This section makes heavy use of the following works [[RC05](#), [RR04](#)].

#### 2.1.1 Motivation and definition

A central goal in Bayesian inference is to evaluate integrals of the type:

$$\mathcal{I} = \int f(x)\pi(x)dx$$

for a given function  $f$ , where  $\pi$  is the posterior. Using the law of large numbers, it would

suffice to have a large sample  $(\{x^{(i)}\}_i)$  from  $\pi$  and use the estimator  $S_n(f) = \frac{1}{n} \sum_i f(x^{(i)})$ . In most instances, however, it is extremely inefficient (or downright impossible in cases where  $\pi$  is only known up to a proportion) to try a direct sampling approach, e.g. with accept/reject type algorithms. Instead we follow a different strategy: obtain an approximate sample from  $\pi$  without directly sampling from it. Markov Chain Monte Carlo methods allow us to do exactly this by using an ergodic Markov chain with stationary distribution  $\pi$ . The underlying Markov chain dynamics means that the samples will be dependent, however the ergodic theorem justifies their use in  $S_n(\cdot)$  just as if they were i.i.d. samples.

Metropolis-Hastings algorithm is one way of constructing such Markov chains. This algorithm makes use of a proposal kernel  $Q$  to construct a kernel  $K$  that has the required stationary distribution  $\pi$ . The outline of the algorithm is as follows; given current state  $x^{(t)}$ ,

1. Generate  $Y_t \sim Q(y|x^{(t)})$

2. Take

$$x^{(t+1)} = \begin{cases} Y_t, & \text{with probability } \rho(x^{(t)}, Y_t) \\ x^{(t)}, & \text{with probability } 1 - \rho(x^{(t)}, Y_t) \end{cases}$$

where

$$\rho(x, y) = \min \left( \frac{\pi(y)Q(y, x)}{\pi(x)Q(x, y)}, 1 \right). \quad (2.1)$$

Intuitively, this algorithm balances the flow out of a state  $x$  into  $y$  with the flow into  $x$  from  $y$ , when the system is in its stationary regime. This balance condition is known as the detailed balance condition and it is one way to generate a reversible chain. On the other hand, the fact that we only need to compute ratios of densities implies that the normalization constant need not be known. We will look at these ideas in more detail in the next subsection.

### 2.1.2 Basic theory

A Markov Chain Monte Carlo method is essentially a Markov Chain specially constructed to have a given limiting distribution. In this subsection, we will discuss the basic properties of Markov Chains and give some basic limit results. Throughout this section, as well as in many other parts of this thesis, we will not explicitly state the  $\sigma$ -algebra of the underlying

probability space to avoid unnecessary cluttering. We will either use the notation  $A$  to denote an event in this  $\sigma$ -algebra (e.g.  $\pi(A)$  for the measure of set  $A$  w.r.t. the measure  $\pi$ ), or use  $\pi(dx)$  when we assume the measure  $\pi$  has a density with respect to some other measure.

### Definition

A *transition kernel* is a function  $K$  such that,

- (i)  $\forall x, K(x, \cdot)$  is a probability measure
- (ii)  $\forall A, K(\cdot, A)$  is measurable.

$n$ -transition kernels can be obtained from this one-step kernel by,

$$K^n(x, A) = \int K^{n-1}(y, A)K(x, dy).$$

Given a transition kernel  $K$ , a sequence  $(X_i)_i$  of random variables is a *Markov chain*, if for any  $t$ , we have;

$$P(X_{t+1} \in A | x_0, \dots, x_t) = P(X_{t+1} \in A | x_t) = \int_A K(x_t, dx).$$

### Basic properties

A Markov chain with a kernel  $K$  is called *irreducible*, if for every  $x$  and  $A$ , there exists  $n$  such that  $K^n(x, A) > 0$ . This property makes sure that any event is reachable from any starting point. One can think of this as a first measure of sensitivity of the Markov chain to its initial conditions.

A set  $C$  is small if there exists  $m$  and a non-zero measure  $\nu_m$  such that,

$$K^m(x, A) \geq \nu_m(A) \quad ; \forall x \in C, \forall A.$$

An irreducible chain has a cycle of length  $d$  if there exists a small set  $C$ , an associated integer

$M$  and a probability distribution  $\nu_M$  such that  $d$  is the greatest common divisor of,

$$\{m \geq 1; \exists \delta_m > 0 \text{ such that } C \text{ is small for } \nu_m \geq \delta_m \nu_M\}.$$

If the chain has a cycle of length 1, it is called *aperiodic*.

Define  $\eta_A = \sum_i I_A(X_i)$ , the number of passages of the Markov chain in  $A$ . A set  $A$  is *Harris recurrent* if  $P_x(\eta_A = \infty) = 1$  for all  $x \in A$ , where  $P_x$  denotes the probability measure of the chain starting at state  $x$ . The chain is Harris recurrent if it is irreducible and there exists a measure  $\psi$  such that for every set  $A$  with  $\psi(A) > 0$ ,  $A$  is Harris recurrent.

A  $\sigma$ -finite measure  $\pi$  is invariant for kernel  $K$  (and the associated chain) if

$$\pi(A) = \int K(x, A) \pi(dx) \quad ; \forall A.$$

If  $\pi$  is a probability measure and the chain is  $\pi$ -invariant, we call the chain *stationary* and  $\pi$  its *stationary distribution*. A stationary Markov chain is *reversible* if it satisfies a condition called the *detailed balance condition* with respect to a density  $\pi$ , i.e.

$$K(y, dx) \pi(dy) = K(x, dy) \pi(dx) \quad ; \forall (x, y).$$

In fact, it is easy to show (see [RC05]) that if the detailed balance condition holds with a probability density  $\pi$ , then  $\pi$  is the invariant density of the chain.

It is important to note that, if a chain has stationary distribution  $\pi$ , it may still fail to converge to stationarity. We give the following example from [RR04]. Suppose the state space is  $X = \{1, 2, 3\}$  with  $\pi(1) = \pi(2) = \pi(3) = 1/3$ . Let  $K(1, \{1\}) = K(1, \{2\}) = K(2, \{2\}) = K(2, \{1\}) = 1/2$  and  $K(3, \{3\}) = 1$ .  $\pi$  is stationary for this kernel, however if the chain starts at state 3, then it cannot transition to other states, so that  $P(X_n = 1) = P(X_n = 2) = 0$  for all  $n > 1$ . This lack of convergence stems from the fact that the states do not communicate with each other, such that the chain gets stuck in a subspace. Irreducibility of a chain, defined earlier, implies that all states properly communicate with one another, eliminating

problematic cases as in the last example. In fact, this concept is central in many convergence theorems.

### Convergence results in fixed dimensions

The basic properties described in the previous section are enough to develop initial convergence results. We stress that these convergence results are “fixed dimension” results, so they are not enough to justify applications in high dimensional inverse problems. However, we will still show them for completeness. High dimensional case will be discussed in next subsection. The first result will use the Total-Variation (TV) distance; we recall its definition <sup>1</sup>:

$$\|\pi_1 - \pi_2\|_{TV} = \sup_A |\pi_1(A) - \pi_2(A)|.$$

**Theorem 10.** *[RC05] Assume the Markov chain has  $\pi$  as its stationary density. If this chain is Harris recurrent and aperiodic, then*

$$\lim_{n \rightarrow \infty} \left\| \int K^n(x, \cdot) \mu(dx) - \pi \right\|_{TV} = 0$$

for every initial distribution  $\mu$ .

This theorem forms the initial justification to use a Markov chain to approximate a (presumably difficult-to-sample) distribution. This result can be made stronger, such that the rate of convergence is uniformly bounded with respect to the initial state. To state this result, we first define the “time of first visit” variable  $\tau_A$ :

$$\tau_A = \inf \{n \geq 1; X_n \in A\}.$$

**Theorem 11.** *[RC05] Assume the Markov chain has  $\pi$  as its stationary density. If the Markov chain is aperiodic, with a small set  $C$  and there exists a real  $\kappa > 1$  such that*

$$\sup_x E_x(\kappa^{\tau_C}) < \infty$$

---

<sup>1</sup>Recall that  $A$  denotes an arbitrary event of the underlying probability space

then,

$$\lim_{n \rightarrow \infty} \sup_x \|K^n(x, \cdot) - \pi\|_{TV} = 0.$$

Recall that  $E_x$  is the expectation with respect to the Markov chain starting at state  $x$ . An equivalent way to define uniformly ergodic chains is the following condition,

$$\|K^n(x, \cdot) - \pi\|_{TV} \leq M\rho^n.$$

Note that, in the above,  $M$  is constant with respect to the initial state  $x$ . This formulation is useful to get a sense of the rate of convergence, independent of the initial state. However, this latter property of independence proves too strong for some chains of practical interest. A weaker formulation corresponds to *geometric ergodicity*. A Markov chain is geometrically ergodic if,

$$\|K^n(x, \cdot) - \pi\|_{TV} \leq M(x)\rho^n,$$

for some  $0 < \rho < 1$ . Essentially, we let the constant  $M$  to depend on the initial state now. Before stating the theorem guaranteeing geometric ergodicity, we need the following definition. A Markov chain is said to satisfy the *drift condition* if there are constants  $0 < \lambda < 1$  and  $b < \infty$ , and a function  $V : X \rightarrow [0, 1]$ , such that,

$$KV \leq \lambda V + b1_C,$$

where  $C$  is a small set. We now proceed to the theorem.

**Theorem 12.** [RR04] *Consider an irreducible, aperiodic Markov chain with a small set  $C$ . Suppose further that the drift condition is satisfied. Then the chain is geometrically ergodic.*

We will end this subsection with a basic central limit theorem. Assume  $g \in L^2(\pi)$ , define  $\hat{g}(X) = g(X) - E_\pi(g(X))$ . We have the following result,

**Theorem 13.** [RC05] *If the Markov chain is aperiodic, irreducible and reversible with invariant distribution  $\pi$ , central limit theorem:*

$$\frac{1}{\sqrt{N}} \sum_i \hat{g}(X_n) \rightarrow_L N(0, \gamma_g^2)$$

applies if,

$$0 < \gamma_g^2 = E_\pi(\hat{g}^2(X_0)) + 2 \sum_i E_\pi(\hat{g}(X_0)\hat{g}(X_k)) < \infty,$$

where  $\rightarrow_L$  denotes convergence in law.

### 2.1.3 MCMC in functional spaces

Our discussion so far did not make explicit reference to the underlying state space of the Markov chain. We tried to keep the overview as general as possible, however this does not imply that the methodology discussed so far applies in the most general case, in particular in the case of infinite dimensional (functional) spaces. The essential problem is that, the Metropolis-Hastings ratio cannot be defined as in Eq. 2.1 anymore, since proposal and target measures need not have a density with respect to Lebesgue measure. In this brief subsection, we will define MCMC algorithm in a more general sense. This will enable us to come up with a MCMC algorithm that works in functional spaces [CRSW13].

Let us define  $\nu(du, dv) = Q(u, dv)\pi(du)$  where  $Q$  is the proposal kernel and  $\pi$  is the target measure (not density!). Similarly, we define  $\nu^*(du, dv) = Q(v, du)\pi(dv)$ . Assuming  $\nu^* \ll \nu$ , we define the Metropolis-Hastings ratio as,

$$\rho(u, v) = \min \left\{ 1, \frac{d\nu^*}{d\nu}(u, v) \right\}.$$

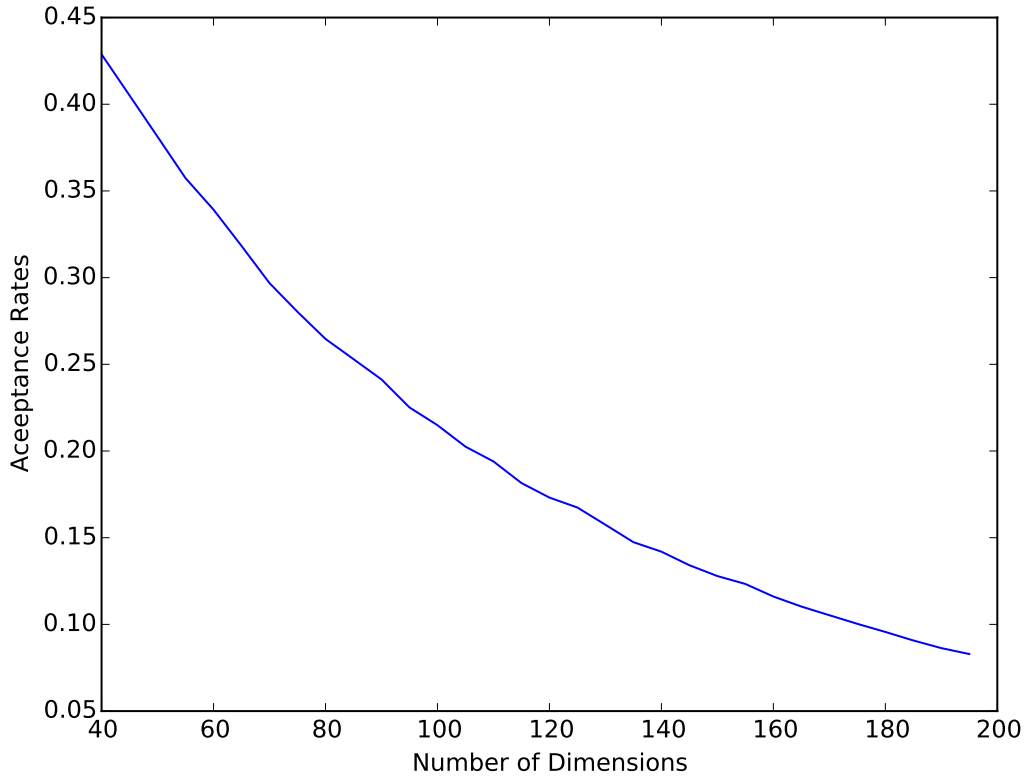
We can easily see that if  $\nu$  and  $\nu^*$  are absolutely continuous with respect to the Lebesgue measure, we recover the previous definition in Eq. 2.1. To see what can go wrong in the functional space setting, we consider the simple random walk proposal MCMC. Assume  $\pi_0$  is  $N(0, C)$ , the proposal reads:

$$v = u + \sqrt{2\delta}\xi$$

with  $\xi \sim N(0, C)$ , i.e. the same distribution as the prior. Here,  $\delta$  is a scaling parameter that can be tuned to optimize the search behaviour of the algorithm. We immediately see a problem: if the target measure is the same as the prior measure, the distribution of  $v$  (the proposed move) when the algorithm is in stationary regime will be  $\sqrt{1 + 2\delta}N(0, C)$  and



therefore the acceptance rates will drop rapidly to 0 as the dimensionality increases; see Figure 2.1 for an example with  $C = I$  and  $\delta = 0.25$ . In fact, the Metropolis-Hastings ratio no longer exists in the infinite dimensional setting, as  $\nu^* \ll \nu$  is not satisfied (i.e. the measures are singular to each other) for this proposal.



**Figure 2.1:** Acceptance rates for the symmetric random walk proposal, when the target is the same as the prior

Let us now consider an alternative proposal:

$$v = (1 - 2\delta)^{1/2}u + \sqrt{2\delta}\xi,$$

with  $\xi \sim N(0, C)$  as before. Even though this proposal looks less natural at first look, it actually leads to a well defined MCMC algorithm with the MH ratio:

$$\rho(u, v) = \min\{1, \exp(\Phi(u) - \Phi(v))\}.$$

Recall that we have  $\frac{d\pi}{d\pi_0}(x) \propto \exp(\Phi(x))$  as in the first chapter. We also see that this proposal preserves the prior and is prior-reversible, i.e. if  $u \sim N(0, C)$  then  $v \sim N(0, C)$ ,  $Q(u, dv)\pi_0(du) = Q(v, du)\pi_0(dv)$ . If the target is the same as prior then the acceptance rate will no longer drop as before (in fact, it will be 1, the reader might compare this to the acceptance rates in Figure 2.1). Such proposals that lead to well defined MH algorithms even in the infinite-dimensional functional setting are called dimensionality-robust. In principle, they should be preferred when the dimensionality of the application is very high, or unknown. However, in practice, the performance of such proposals depend on how well the posterior measure is approximated by the prior. Hence the theoretical safety that such an approach brings does not necessarily translate to universal applicability. It is also interesting to note that when  $\delta = 1/2$  the algorithm becomes an independence sampler. This will be relevant when we discuss multi-resolution MCMC strategies in functional setting in the last chapter.

#### 2.1.4 Tempered Monte Carlo methods

We end our brief review of MCMC methodology by tying it to particle-based methods. The basic idea of using the plain Metropolis-Hastings algorithm directly on the density of interest only works if the said density is simple enough. Densities with multiple modes or highly peaked regions need to be handled differently, as even when the algorithm has a guarantee to converge, the convergence can be extremely slow. It is very usual to observe a MCMC method to stay stuck in a region of space for a very long time in such difficult problems. The two methods that we will discuss in this subsection attempt to solve this issue. In each case, the new goal is to simulate from a sequence of non-normalized densities  $(\pi_i)_{i=1}^m$  on the same state space. The index  $i$  is called the “temperature”,  $\pi_1$  being the “cold” density while  $\pi_m$  is the hot one. These densities are modifications of the target density of interest,  $\pi$ . We choose them such that the cold density starts very flat and gradually heats up to the target density  $\pi_m = \pi$ . The motivation is to utilize the flatness of the cold densities to enable faster exploration of the state space, and guide the system towards high density regions of the original density of interest. We will discuss a particular way of constructing such a sequence in the next chapter.

The first method we discuss is called Simulated Tempering [GT95]. In this method, we consider the augmented state space  $(x, i) \in \mathcal{X} \times \{0, 1, \dots, m\}$ , where the temperature is now taken as random. The stationary distribution of the sampler is proportional to  $\pi_i(x)h(i)$ , where  $h(i)$  is called the “pseudo-prior”. One iteration of this algorithm looks as follows [GT95]:

1. Update  $x$  using Metropolis-Hastings for  $\pi_i(\cdot)$ .
2. Set the proposed new temperature  $j = i \pm 1$  according to probabilities  $q_{i,j}$  where  $q_{1,2} = q_{m,m-1} = 1$  and for other indices  $q_{i,j} = \frac{1}{2}$ .
3. Calculate the MH ratio,

$$\rho = \frac{\pi_j(x)h(j)q_{j,i}}{\pi_i(x)h(i)q_{i,j}},$$

and accept/reject transition as per the usual Metropolis-Hastings method.

When the tempered density sequence and the pseudo-prior are chosen well, this algorithm can converge very quickly, where a standard MCMC would fail. It turns out, however the choice of these are far from trivial. The paper by Geyer et.al. [GT95] proposes a rather contrived mechanism to adaptively find these parameters. Parallel tempering MCMC, to be discussed next, alleviates part of this problem by not requiring a pseudo-prior. Though choice of a good sequence of densities remains an issue still. We will see in the next chapter that such a choice becomes much easier in the context of Sequential Monte Carlo.

Population based MCMC [JSH07] replaces the single chain approach of simulated tempering with multiple interacting chains that run in parallel. Instead of sampling from the augmented space of  $(X, i)$ , we now try to sample from the following target measure:

$$\pi^*(x_{1:m})dx_{1:m} = \Pi_n \pi_n(x_n)dx_{1:m}$$

where again  $\pi_1$  is the cold density that is flat and  $\pi_m = \pi$  is the hot density with complicated features. Our goal, then, is to construct a Markov kernel that is  $\pi^*$ -irreducible, aperiodic and admits  $\pi^*$  as its stationary distribution. We achieve this by considering the target as having vector components  $(x_1, \dots, x_m)$ . We visualize the resulting algorithm as running parallel MCMC for the different components, while exchanging information between the chains

at certain points. Once the proposed move kernel is fixed, Metropolis-Hastings algorithm is applied as discussed previously. We now give some example move proposals from the literature;

1. *Mutation*: This move seeks to update a single component of the target using a Markov kernel, i.e.  $X_{n+1}|x_n \sim K(x_n, \cdot)$ , which is then applied independently to all components.
2. *Exchange*: This move attempts to change the values between two randomly selected chains  $i$  and  $j$ .
3. *Crossover*: It is possible to attempt a less extreme swap than the Exchange kernel, if each component  $x_i$  itself is a  $d$ -dimensional vector  $x_i = (x_{i1}, \dots, x_{id})$ . This move attempts to crossover the  $l$ th position in the vector for chains  $i$  and  $j$  such that the proposed change would look like  $x_{n+1,i} = (x_{n,i1}, \dots, x_{n,jl}, \dots, x_{n,id})$  and  $x_{n+1,j} = (x_{n,j1}, \dots, x_{n,il}, \dots, x_{n,jd})$ .

The algorithm would typically choose and mix some of these possible moves, to get a composite move at each time step. It is known that the simulated tempering algorithm may converge faster, if an efficient way to come up with a good pseudoprior is available. This now ties the discussion to Monte Carlo methods based on interacting particle systems, otherwise known as Sequential Monte Carlo methods.

## 2.2 Sequential Monte Carlo Samplers

The most popular sampling method by far is the Markov Chain Monte Carlo (MCMC) method. A less well known competition comes in the form of Sequential Monte Carlo (SMC). One of the goals of this thesis is to investigate the practical properties of SMC, especially in the high dimensional setting and try to answer this question: “Can SMC be made practically relevant in high dimensional inverse problems?”. This section will review some basic material about SMC, including basic theory using Feynman-Kac path measures. The discussion is based on the paper by Del Moral et. al. [DMDJ06] and the seminal book by the same author [DM04].

### 2.2.1 Motivation and definition

The goal is to sample from a sequence of densities  $(\pi_1, \dots, \pi_n)$ . In a Bayesian setup, this sequence may correspond to increasing number of observations, different levels of annealing, etc. For example, the former case is related to radar tracking of moving objects, which is a canonical problem in the signal processing area.

#### Hidden Markov Models

The canonical example to motivate Sequential Monte Carlo methods (or “particle filters” as they are known in applied fields) is Bayesian inference in Hidden Markov Models [DJ09]. A Hidden Markov Model (HMM)  $\{X_n, Y_n\}_n$  is made up of a hidden Markovian dynamics  $(X_n | X_{n-1} = x_{n-1} \sim \pi_{X_n | X_{n-1}}(\cdot | x_{n-1}))$  and a conditionally independent observation process  $(Y_n | X_n \sim \pi_{Y | X_n}(\cdot | x_n))$ . The goal is to make inference on the hidden variables  $\{X_n\}_n$  conditioned upon the observations  $\{Y_n\}_n$ . When the transition dynamics are linear (i.e.  $X_n = AX_{n-1} + BV_n$  for matrices  $A, B$  and  $V_n \sim N(0, I)$ ) and the observation process is linear ( $Y_n = CX_n + DW_n$ , where  $C, D$  are some matrices and  $W_n \sim N(0, I)$ ), the maximum likelihood inference problem can be solved by Kalman filters. The Extended Kalman Filters attempts to solve the general problem using linearization. Particle filters, on the other hand, enables inference without such possibly crude approximations. Put more concretely, the inference goal is to extract information from the following sequence of densities:

$$\pi_n(x_{1:n} | y_{1:n}) \propto \pi_1(x_1) \prod_{i=2}^n \pi_{X_i | X_{i-1}}(x_i | x_{i-1}) \pi_{Y | X_i}(y_i | x_i).$$

#### SMC Samplers

The SMC way of sampling from a sequence of target densities consists of having a weighted particle system at each time index, i.e.  $\{x_n^i, w_n^i\}_n^N$ , with  $N$  being the number of particles. In this system, at time  $n$ ,  $\{x_n^i\}$  are distributed approximately according to  $\pi_{n-1}$ , i.e.  $\sum_i \delta_{x_n^i} \approx \pi_{n-1}$ ,  $\sum_i w_n^i \delta_{x_n^i} \approx \pi_n$  with  $\sum_i w_n^i = 1$ . The last equality indicates that the weights are normalized to sum to 1. For this reason, in our discussion we will only be concerned to find

them up to a proportionality constant (i.e. the non-normalized  $\omega_n^i \propto w_n^i$ ) which is the inverse of the total non-normalized weights (i.e.  $w_n^i = \frac{\omega_n^i}{\sum_i \omega_n^i}$ ). We assume that the first target density is easy to sample from and that we have a reasonable prediction kernel  $K_n$ , so that  $\pi_{n-1}K_n(A) = \int K_n(x_{n-1}, A)\pi_{n-1}(x_{n-1})dx_{n-1} \approx \pi_n(A)$ . To motivate the SMC procedure, let us first discuss the basic importance sampling solution. Using the aforementioned assumptions, the importance distribution at time  $n$  is  $\nu_n = \pi_1 \prod_{j=1}^n K_j$  and hence the (non-normalized) weights can be calculated as  $\omega = \frac{\pi_n}{\nu_n} = \frac{\pi_n}{\pi_1 \prod_{j=1}^n K_j}$ . If the importance distribution could be computed pointwise efficiently, then this would be a reasonable procedure. That is not the case in the majority of applications and hence we need an alternative approach. That alternative is to consider an artificial joint density, for which  $(\pi_n)_n$  are marginals, so that a simpler recursive formula for the weights can be derived.

In order to create the required joint distribution, we employ backward kernels  $L_k$ , so that the joint density at time  $n$  reads  $\tilde{\pi}_n(x_{1:n}) = \pi_n(x_n) \prod_{k=1}^{n-1} L_k(x_{k+1}, x_k)$ . It can easily be seen that such a construction admits  $\pi_n$  as marginal. Under this scheme, the weights can be expressed as;

$$\omega_n(x_{1:n}) = \frac{\tilde{\pi}_n(x_{1:n})}{\nu_n(x_{1:n})} = \omega_{n-1}(x_{1:n-1}) \tilde{\omega}_n(x_{n-1}, x_n)$$

with  $\nu_n(x_{1:n}) = \pi_1(x_1) \prod_k K_k(x_{k-1}, x_k)$  and  $\tilde{\omega}_n(x_{n-1}, x_n) = \frac{L_{n-1}(x_n, x_{n-1})\pi_n(x_n)}{K_n(x_{n-1}, x_n)\pi_n(x_{n-1})}$ . Essentially, we get rid of the marginalization step of  $\pi_{n-1}K_n = \int x_{n-1}K_n(x_{n-1}, \cdot)\pi_{n-1}(x_{n-1})dx_{n-1}$  and replace it with the products above.

The next step, then, is to figure out the best, in some sense, backward kernels. A popular criterion in importance sampling literature is the variance of weights, and it is reasonable to select the kernels that minimize this. We will now briefly follow the discussion in Section 3.3 of Del Moral et.al. [DMDJ06]. In the development of SMC outlined above, we essentially replace IS in 1-dimension with IS in  $n$ -dimensions. This is used to alleviate the problem of computing the marginal of proposal density. On the other hand, such an enlargement of space leads to increased variance of importance weights. Hence, unsurprisingly, the optimal kernel is the one that takes us back to the 1-dimensional case, i.e. ;

$$L_{k-1}^{opt}(x_k, x_{k-1}) = \frac{\nu_{k-1}(x_{k-1})K_k(x_{k-1}, x_k)}{\nu_k(x_k)}$$

which is impractical for the same reasons as IS (i.e. the difficulty of exactly calculating  $\nu_k(x_k)$ ); but it gives us a way of constructing an approximately optimal kernel. As a first step, we replace  $\nu_{k-1}$  with  $\pi_k$ , the argument being that they should be close to each other if the importance density is any good. In this case the kernel reads;

$$L_{k-1}^{opt'}(x_k, x_{k-1}) = \frac{\pi_{k-1}(x_{k-1})K_k(x_{k-1}, x_k)}{\pi_{k-1}K_k(x_k)}.$$

In the final step, we assume  $K_n$  is  $\pi_n$  invariant, which is usually the case in practice, therefore  $\pi_{k-1}K_k(x_k) \approx \pi_k$ . The final expression reads;

$$L_{k-1}^{opt''}(x_k, x_{k-1}) = \frac{\pi_{k-1}(x_{k-1})K_k(x_{k-1}, x_k)}{\pi_k(x_k)}$$

where the associated incremental weights are simply;

$$\tilde{\omega}_n(x_{n-1}, x_n) = \frac{\gamma_n(x_{n-1})}{\gamma_{n-1}(x_{n-1})}$$

with  $\pi_n \propto \gamma_n$ . In the case where  $\pi_n$  and  $\pi_{n-1}$  are posteriors with a common prior,  $\gamma_n$  is simply the likelihood at time  $n$ .

Let us end this discussion with the structure of a typical implementation of the ideas presented in this section. Assume that at time  $t-1$ , we have a particle system  $\{x_{t-1}^i, \omega_{t-1}^i\}_{i=1}^{N_{t-1}}$  with uniform weights (i.e.  $\omega_{t-1}^i = 1$  and  $w_{t-1}^i = \frac{1}{N_{t-1}}$ , with  $N_{t-1}$  the number of particles at time  $t-1$ ), the new particle system is generated as follows,

1. *Correction*: Assign new weights to the particles by,

$$\omega_t^i = \tilde{\omega}_t(x_{t-1}, x_t) = \frac{\gamma_t(x_{t-1})}{\gamma_{t-1}(x_{t-1})}.$$

2. *Selection*: Resample, according to multinomial sampling<sup>2</sup> (i.e.  $\hat{x}_t \sim \text{Multinomial}((\tilde{w}_t^i))_i$ ), using the normalized weights ( $\tilde{w}_t^i = \frac{\tilde{\omega}_t^i}{\sum_i \tilde{\omega}_t^i}$ ) as selection probabilities. Then re-set the weights to 1.

---

<sup>2</sup>Other resampling methods are also available.

3. *Mutation*: For each particle  $\hat{x}_{t-1}$ , sample new particles  $x_t$  from  $K_n$ , where we assume  $K_n$  is  $\pi_n$  invariant.

Figure 3.2 in the next chapter gives an example run of this algorithm. Note that the target density sequence is adaptively determined by the algorithm in that example.

### 2.2.2 Curse of Dimensionality Discussions

Sequential Monte Carlo Samplers have been utilized in many different domains. A frequently observed problem is the “weight degeneracy” or the “collapse” of the particle system, where only a handful particles have significant weights at the later iterations. Many remedies are proposed in the literature, including resampling and tempering which are discussed in the next chapter. Bickel et.al. [BLB08], however, showed that in a simple Gaussian setup, the degeneracy phenomenon is inherent in the algorithm (in particular, inherent to the reweighing step), when the dimensionality is large. In this subsection, we briefly discuss their results and a recent paper taking on a different perspective [RvH15]. We will see that these two viewpoints converge on the idea that what really matters is the “effective dimension” of the system. This result is crucial in understanding the effectiveness of the proposed algorithms in this thesis.

In the aforementioned work, Bickel et.al. [BLB08] examine the behaviour of the importance weights as the system dimension and the sample size increases. Their results imply that, to avoid collapse, the sample size must grow super-exponentially in the effective dimension. They also conjecture that methods such as deterministic simulated tempering (as opposed to our adaptive tempering in the next chapter) does not provide a remedy.

Let us briefly overview their setup and discuss the results. The data model they use is a linear Gaussian-Gaussian model, i.e.  $Y = HX + \varepsilon$ , where  $Y$  is a  $d \times 1$  observation vector,  $H$  is a known  $d \times q$  matrix and  $X$  is the  $q \times 1$  parameter vector. Both the prior and the error distribution are Gaussian, with  $X \sim N(\mu_X, \Sigma_X)$  and  $\varepsilon \sim N(0, \Sigma_\varepsilon)$  and the noise is independent of the parameter. Finally, let  $d' = \text{rank}(H)$  and  $(\lambda_i^2)_i$  be the singular values of  $\text{cov}(HX)$ . We have the following result [BLB08].

**Theorem 14.** *Assume the following,*



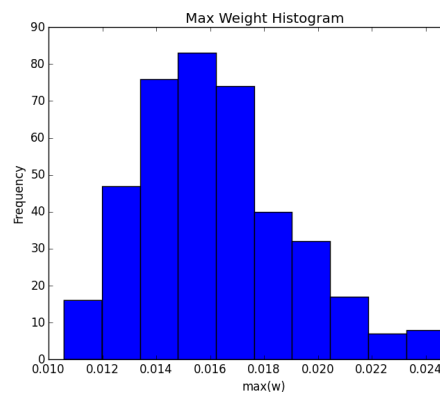
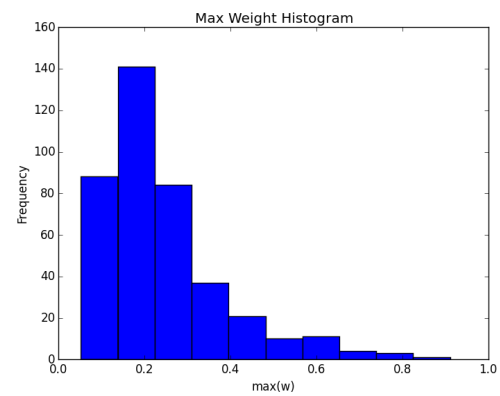
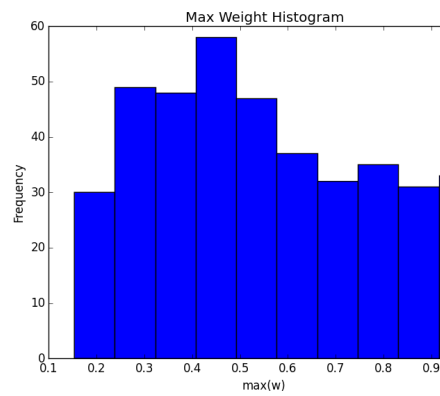
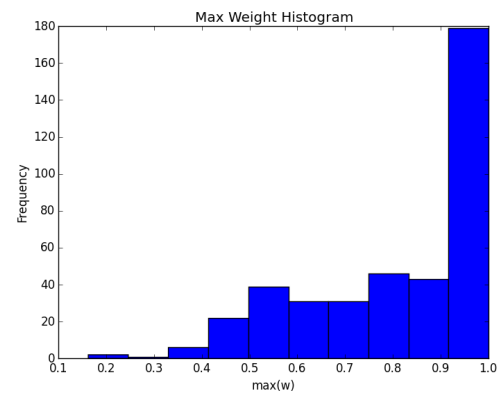
A1. There is a positive constant  $\delta$  s.t.  $\frac{1}{\delta} \geq \lambda_1 \geq \lambda_2 \geq \dots \geq \delta$ ,

A2.  $\frac{2}{d'} \sum_{j=1}^{d'} (3\lambda_j^4 + 2\lambda_j^2) \rightarrow \sigma^2 > 0$ .

Then, if  $\frac{\log n \log d'}{d'} \rightarrow 0$ , we have the largest normalized weight  $\omega(n) \rightarrow 1$  in probability.

This result implies that unless we have exponential growth of the sample size in terms of  $d'$ , we have weight collapse, i.e. the maximum weight converges to 1. The first assumption is very important, as it essentially means “when the effective dimensionality is high”. Bickel et.al. [BBL08] take this further and show that when  $\sum_i \lambda_i^2 < \infty$ , there is no weight collapse. They define this sum as the effective dimension of  $X$ . The reason for weight collapse is that the proposal distribution and the target density become mutually singular [BBL08]. We will now replicate a small experiment from Bengtsson [BBL08] that showcases the weight degeneracy behaviour of a simple importance sampler in increasing dimensions. We assume a simple linear problem with Gaussian noise, i.e.  $Y = X + \varepsilon$  such that  $X \in \mathbb{R}^d$ ,  $\varepsilon \sim N(0, I)$  and we assume  $n$  observations. If we use the prior as the proposal distribution, the importance weights become  $w_i \propto \exp(-\frac{1}{2} \sum_{j=1}^d \varepsilon_j^2)$ . Hence, we can sample the maximum weight in an importance sampler by first generating  $n$  weights as above and taking their maximum. We repeat this exercise for 400 times with  $n = 1000$  and various  $d$  and the results are shown in Figure 2.6. We observe that at  $d = 100$  we start to observe few particles dominating the system, and at  $d = 400$  most of the simulations essentially collapse with the maximum weight very close to 1.

This “curse of dimensionality” does not effect the classical tracking/filtering problems, which are often not high dimensional. On the other hand, Bickel’s result caused a lot of pessimism in the viability of SMC in high dimensional problems. A more optimist discussion, together with proof-of-concept theory, can be found in Rebeschini et.al. [RvH15]. In this paper, they argue that it is, in principle, possible to develop local particle filters whose local approximation error is dimension-free. Another way of putting it is that, even though the overall error explodes as the effective dimension increases, the local error (i.e. error w.r.t. to a marginal) may not. A simple though experiment shows this idea: imagine a high dimensional hidden Markov model (HMM). Sequential Monte Carlo method, also called Particle Filters, is a natural inference tool for this model. Now assume that, in this HMM all the dimensions

Figure 2.2:  $d = 10$ Figure 2.3:  $d = 40$ Figure 2.4:  $d = 100$ Figure 2.5:  $d = 400$ Figure 2.6:  $\max(w)$  plots for different dimensions  $d$ .

are i.i.d., in other word this high dimensional chain is actually made up of unrelated and identical one-dimensional HMMs. We observe that, since the effective dimensionality is exactly the same as the actual dimensionality, Bickel's result applies. On the other hand, if we take any one of these chains, i.e. considering marginals on single parameters, and perform SMC inference on these individual chains, the local error is independent of the overall dimensionality. They also demonstrate that this result still holds when the complete independence condition is weakened to a decay of correlation among parameters. We refer to [RvH15] for details.

In sum, we have seen that the error explodes exponentially as the effective dimensionality increases. But it is still possible for the effective dimensionality to be bounded while the actual dimensions increase. Such a case is possible when observations effect only a small amount of parameters significantly, while most of the parameters are highly influenced by the prior. Finally, even when the effective dimensionality is very high, it may be possible to have dimension independent error rates for the marginals. These results motivate us to further investigate the practical performance of SMC in inverse problems.

### 2.2.3 A basic Central Limit Theorem for SMC

In this subsection, we will briefly state and discuss the central limit theorem for sequential Monte Carlo methods by Chopin [Cho04]. We begin by recursively defining the following quantities, beginning with  $\tilde{V}_0(g) = \text{Var}_{\pi_0}(g)$ , with  $g$  being any measurable function,

$$V_t(g) = \tilde{V}_t(\tilde{\omega}_t \cdot (g - E_{\pi_t}(g))),$$

$$\hat{V}_t(g) = V_t(g) + \text{Var}_{\pi_t}(g),$$

$$\tilde{V}_t(g) = \hat{V}_t(E_{K_t}(g)) + E_{\pi_t} \text{Var}_{K_t}(g).$$

Note that these quantities need not be finite for any  $t$  in general. We now define the space of functions  $g$  for which the central limit theorem will be stated. Let us define  $\Phi_t$  recursively

to be the set of measurable functions such that for some  $\delta > 0$ ,

$$E_{\pi_{t-1}} \|\tilde{\omega}_t \cdot g\|^{2+\delta} < \infty$$

and that the function  $E_{K_t(x_{t-1}, \cdot)}(\tilde{\omega}_t g(\cdot))$  is in  $\Phi_{t-1}$ . The initial set  $\Phi_0$  contains all measurable functions with finite second moments with respect to  $\pi_0$ . We can now present the following theorem [Cho04]:

**Theorem 15.** *If the selection step <sup>3</sup> is Multinomial resampling, and if the unit function belongs to  $\Phi_t$  for every  $t$ , then for any  $g \in \Phi_t$ ,  $E_{\pi_t}(g)$ ,  $V_t(g)$  and  $\hat{V}_t(g)$  are finite quantities and the following convergence results hold as  $N \rightarrow \infty$ :*

$$N^{1/2} \left\{ \frac{\sum_i \tilde{\omega}_t^i g(x_t^i)}{\sum_i \tilde{\omega}_t^i} - E_{\pi_t}(g) \right\} \rightarrow_L N(0, V_t(g))$$

and,

$$N^{1/2} \left\{ \frac{\sum_i g(\hat{x}_t^i)}{N} - E_{\pi_t}(g) \right\} \rightarrow_L N(0, \hat{V}_t(g)).$$

Before closing this section, it is worth noting that estimating the (asymptotic) variance of an SMC estimator (i.e.  $V_t(g)$  and  $\hat{V}_t(g)$ ) is difficult and it is an active research area, for a recent work on this subject we refer to [LW15].

#### 2.2.4 Feynman-Kac measures and some results

Feynman-Kac measures generalize Markov chains by incorporating “potentials”. In effect, these measures form the basis of a wide variety of statistical models as well as open up new avenues in probability research. The covered areas include broad topics as Bayesian statistics, in addition to signal filtering, genetic/evolutionary algorithms, particle physics and sequential Monte Carlo. In this subsection, we will define the Feynman-Kac prediction flow, McKean interpretation and the interacting particle systems, which will form the theoretical basis for sequential Monte Carlo. We will end this section with a time-uniform convergence theorem which will then be employed in the next chapters. All the material in this section is

---

<sup>3</sup>Recall the discussion at the end of Section 2.2.1.

prepared using the seminal book by Del Moral [DM04], in particular sections 3.1, 3.2 and 7.4.3.

The basic building blocks are a sequence of potential functions  $G_n$  and a Markov chain  $(X_n)_n$ . The sequence of Feynman-Kac prediction measures is defined as follows;

$$\nu_n(f_n) = \frac{\gamma_n(f_n)}{\nu_n(1)} \text{ with } \gamma_n(f_n) = E_{\nu_0}(f_n(X_n) \prod_{p=0}^{n-1} G_p(X_p))$$

where  $\nu_0$  denotes the initial measure, the prior in our case. We assume the potential functions are bounded and positive. Intuitively, the Markov chain corresponds to the predictive kernels  $M_n$  in the SMC algorithms, while the potentials corresponds to weights. In this view, Feynman-Kac measure gives the ideal behaviour of a particle that evolves according to  $M_n$  and is selected/sampled/killed according to  $G_n$ . From the opposite view, interacting particle systems, and SMC in particular, can be seen as a stochastic linearization of this ideal measure. This sequence of prediction measures satisfy the following nonlinear recursive equation;

$$\nu_n = \nu_n K_{n+1, \nu_n}$$

where  $K_{n+1, \nu_n}$  is a nonunique collection of Markov kernels, called the McKean interpretation of the underlying Feynman-Kac measure. This McKean kernel can be decomposed into selection  $(S_{n, \nu})$  and mutation kernels  $M_{n+1}$ , i.e.

$$K_{n+1, \nu} = S_{n, \nu} M_{n+1}.$$

To make it more concrete, let us consider kernels corresponding to multinomial and residual sampling schemes, respectively they read;

$$S_{n, \nu} = \Psi_n(\nu)$$

$$S_{n, \nu} = \epsilon_n G_n(x_n) \delta_{x_n} + (1 - \epsilon_n G_n(x_n)) \Psi_n(\nu)$$

with,

$$\Psi_n(\nu)(dx) = \frac{1}{\nu(G_n)} G_n(x) \nu(dx).$$

We will now define the N-interacting particle system (NIPC) corresponding to McKean kernel  $K_{n,\nu}$ . NIPC is a Markov chain that takes values in the product space  $E^N$ , i.e.  $\xi_n^N = (\xi_n^{N,1}, \dots, \xi_n^{N,N}) \in E^N$  with  $\xi_n^{N,i}$  denoting the  $i$ th particle at  $n$ th epoch. The initial configuration  $\xi_0^N$  consists of  $N$  i.i.d. particles drawn from  $\nu_0$ . The transitions of NIPC is given by,

$$P_{\nu_0}^N(\xi_n^N \in dx_n | \xi_{n-1}^N) = \prod_{p=1}^N K_{n, \frac{1}{N} \sum_{i=1}^N \delta_{\xi_{n-1}^{N,i}}}(\xi_{n-1}^{N,p}, dx_n^p).$$

We see from this that the previous deterministic recursive equations of the form  $\nu_n = \nu_{n-1} K_{n, \nu_{n-1}}$  are now replaced with stochastic  $\xi_n^N = \xi_{n-1}^N K_{n, \frac{1}{N} \sum_{i=1}^N \delta_{\xi_{n-1}^{N,i}}}$ . If we associate  $\xi_n^N$  with the corresponding empirical measure  $\hat{\nu}_n$ , we can also interpret the NIPC as a measure valued Markov chain. We will see that, under certain assumptions, these random measures actually approximate  $\nu_n$ , i.e. the Feynman-Kac measure sequence corresponding to the same selection and mutation kernels. We now briefly mention a related theorem.

**Theorem 16.** *Assume the following;*

1. *There exists a sequence of strictly positive numbers  $\epsilon_n \in (0, 1]$  s.t. for any  $x_n, y_n$ ;*  
 $G_n(x_n) \geq \epsilon_n G_n(y_n) > 0$ .
2. *There exists some integer  $m \geq 1$  and some sequence of numbers  $\epsilon_n \in (0, 1)$  s.t. for any  $p$  and  $x_n, y_n$  we have  $M_{p,p+m}(x_p, \cdot) = M_{p+1} M_{p+2} \dots M_{p+m}(x_p, \cdot) \geq \epsilon_n M_{p,p+1}(y_p, \cdot)$*

*Then we have the following result;*

$$\sup_{n \geq 0} \sup_{f_n \in \text{Osc}_1} \sqrt{N} E(|[\nu_n^N - \nu_n](f_n)|^p)^{\frac{1}{p}} \leq c(p, \epsilon)$$

*where  $\text{Osc}_1$  is the class of functions s.t.  $\max(f) - \min(f) \leq 1$  and  $c(p, \epsilon)$  denotes a constant depending on  $p, \epsilon$  as well as the constants in the aforementioned assumptions.*

This very useful time-uniform convergence result is then used to prove consistency of various SMC algorithms. Time-uniformity means that we do not have to explicitly take into account the number of steps taken in the algorithm. This is very useful when this quantity is unknown beforehand, as in the case of our adaptive SMC algorithm explained in the next chapter.

# CHAPTER 3

## Uni-Resolution Adaptive Sequential Monte Carlo Inference

Inverse problems, as defined in the first chapter, can be attacked from different angles. One essential feature of the problem is the computational complexity of calculating the likelihood; hence efficient allocation of computational resources is imperative. Efficient exploration of high dimensional posteriors constitute the other significant challenge [MWBG12]. This chapter is concerned with an adaptive SMC approach that attempts to address exactly these issues. We will utilize the emerging literature on dimension independent mixing MCMC, in addition to our novel approaches to SMC adaptation. Our contributions to SMC adaptation addresses two issues: first is the common criticism of the large number of parameters that have to be pre-determined. The second issue is the proper allocation of computational resources when sampling from a tempered sequence of densities. Note that in this chapter we consider the forward problem as “black-box”, hence no particular structure of the underlying PDE is exploited here. We will consider an alternative approach, which utilizes the mathematical properties of the underlying forward problem, in the next chapter.

### 3.1 Problem Statement

Discrete approximations (sometimes called “meshes” in case of spatial approximations) are often used in practice to obtain approximate solutions to complicated dynamical systems, e.g. Galerkin projection in the case of PDEs. The required refinement level of the mesh may depend on multitude of practical factors and hence, from a methodological perspective, it is

important to provide inference algorithms that are robust to remeshing (i.e. changing the resolution of the mesh). In our scenario, such a result is proved for MCMC ([HSS12]) and SMC ([BJMS15]) methods. We will briefly state and discuss the SMC result in the following sections.

Another important consideration is the behaviour of the inference method in large-data or small-error regimes. Appropriate tuning of algorithmic parameters, such as the scale of a Gaussian random walk kernel, becomes crucial. Adaptive MCMC methods solve this by adapting kernel parameters using history of the chain. However, the inherent dependency of different steps in a MCMC method makes the task of analyzing and developing such methods very difficult.

A different approach which we will adopt here, involves adaptive SMC samplers [DMDJ06]. In particular, we will consider a sequence of measures which interpolate from prior to posterior where the sequential nature of the approximating particle system allows for smooth evolution of particle distribution and weights from the typically simple prior to the potentially very complex posterior. Recent work in the context of inverse problems (Kantas et.al. [KBJ13]) has shown how, using the aforementioned dimension independent MCMC methods within SMC, it is possible to construct algorithms which combine the dimension independent aspects of novel MCMC algorithms with the desirable self-adaptation of particle methods.

## 3.2 Tempered Sequential Monte Carlo

In the second chapter, we discussed how the particle system of SMC follows a sequence of densities. We now introduce a sequence of “bridging” densities which enable us to connect  $\nu_0$  (prior measure) to  $\nu^y$  (posterior measure). Assuming the Bayes theorem holds with  $\frac{d\nu^y}{d\nu_0}(u) \propto \pi(u)$ , the bridging densities are constructed as follows:

$$\pi_n \propto \pi(u)^{\phi_n}$$

where  $0 = \phi_0 < \phi_1 < \dots < \phi_p = 1$ ; we refer to  $\phi_j$  as temperatures. We let  $\nu_n$  denote the probability measure with density proportional to  $\pi_n$  with respect to  $\nu_0$ . Assuming  $\pi(u)$



is finite  $\nu_0$  almost surely, as is true for the uniform prior setup explained in the previous chapter, we obtain:

$$\frac{d\nu_n}{d\nu_0}(u) \propto \pi(u)^{\phi_n}, \quad \frac{d\nu_n}{d\nu_{n-1}}(u) \propto l_{n-1}(u) := \pi(u)^{\phi_n - \phi_{n-1}}.$$

Although  $\nu = \nu_p$  may be far from  $\nu_0$ , careful choice of the  $\phi_n$  can ensure that  $\nu_n$  is close to  $\nu_{n-1}$  allowing gradual evolution of the particle approximation of  $\nu_0$  to that of  $\nu$ . Other choices of bridging densities are possible and are discussed in e.g. Del Moral et.al. [DMDJ06].

The overview of the algorithm is shown in Figure 3.1.

- 
0. Sample  $\{u_0^m\}_{m=1}^M$  i.i.d. from  $\nu_0$  and define the weights  $w_0^m = M^{-1}$  for  $m = 1, \dots, M$ . Set  $n = 1$  and  $l = 0$ .
  1. For each  $m$  set  $\hat{w}_n^m = \ell_{n-1}(u_{n-1}^m)w_{n-1}^m$  and sample  $u_n^m$  from  $K_n(u_{n-1}^m, \cdot)$ ; calculate the normalized weights

$$w_n^m = \hat{w}_n^m / \left( \sum_{m=1}^M \hat{w}_n^m \right).$$

2. Calculate the Effective Sample Size (ESS):

$$ESS_{(n)}(M) := \frac{\left( \sum_{m=1}^M w_n^m \right)^2}{\sum_{m=1}^M (w_n^m)^2}. \quad (3.1)$$

If  $ESS_{(n)}(M) \leq M_{\text{thres}}$ :

resample  $\{u_n^m\}_{m=1}^M$  according to the normalized weights  $\{w_n^m\}_{m=1}^M$ ;  
 re-initialise the weights by setting  $w_n^m = M^{-1}$  for  $m = 1, \dots, M$ ;  
 let  $\{u_n^m\}_{m=1}^M$  now denote the resampled particles.

3. If  $n < p$  set  $n = n + 1$  and return to Step 1; otherwise stop.
- 

**Figure 3.1:** Standard SMC Samplers.  $M_{\text{thres}} \in \{1, \dots, M\}$  is a user defined parameter.

### 3.3 Convergence Properties of Non-Adaptive SMC

The issue of dimensionality in SMC methods has attracted substantial attention in the literature [BCJ14, BCJW11, RvH15]. A result for our setup, which we will briefly state in this section, can be found in Beskos et.al. [BJMS15].

We assume that there exists  $\kappa > 0$  such that for each  $n \geq 0$  and  $u$ ;

$$\kappa \leq l_n(u) \leq 1/\kappa. \quad (3.2)$$

We note that this holds for the elliptic inverse problem discussed before, when the uniform prior is employed.

Let  $P$  denote the collection of all probability measures on our probability space  $E$ . Let  $\mu = \mu(w)$  and  $\nu = \nu(w)$  denote two possibly random elements in  $P$ . We define the distance between  $\mu, \nu \in P$  by

$$d(\mu, \nu) = \sup_{|f|_\infty \leq 1} \sqrt{E^w |\mu(f) - \nu(f)|^2}.$$

This definition of distance is indeed a metric on the space of random probability measures; in particular it satisfies the triangle inequality. In the context of SMC, the randomness stems from various sampling operations within the algorithm.

Bekos et.al. [BJMS15] provides the following convergence result for non-adaptive SMC;

**Theorem 17.** *Assume 3.2. Consider a non-adaptive SMC which resamples at every iteration. Then, for any  $n \geq 0$ ,*

$$d(\nu_n^M, \nu_n) \leq \sum_j (2\kappa^{-2})^j \frac{1}{\sqrt{M}},$$

where  $\nu_n^M$  denotes the empirical measure corresponding to the  $n$ th iteration of the SMC particle system with  $M$  particles.

We make some comments about this.

- The measure  $\nu_p$  is well approximated by  $\nu_p^M$  in the sense that, as the number of particles  $M \rightarrow \infty$ , the approximating measure converges to the true measure. The result holds in the infinite dimensional setting. As a consequence, the algorithm as stated is robust to finite dimensional approximation.
- In principle, the theory applies even if we skip the move steps. However, moving the particles according to a non-trivial  $\nu_n$ -invariant measure is absolutely essential for the methodology to work in practice. This can be seen by noting that if we skip all the move steps, the final particle system becomes a weighted set of samples from the prior,

clearly undesirable in general.

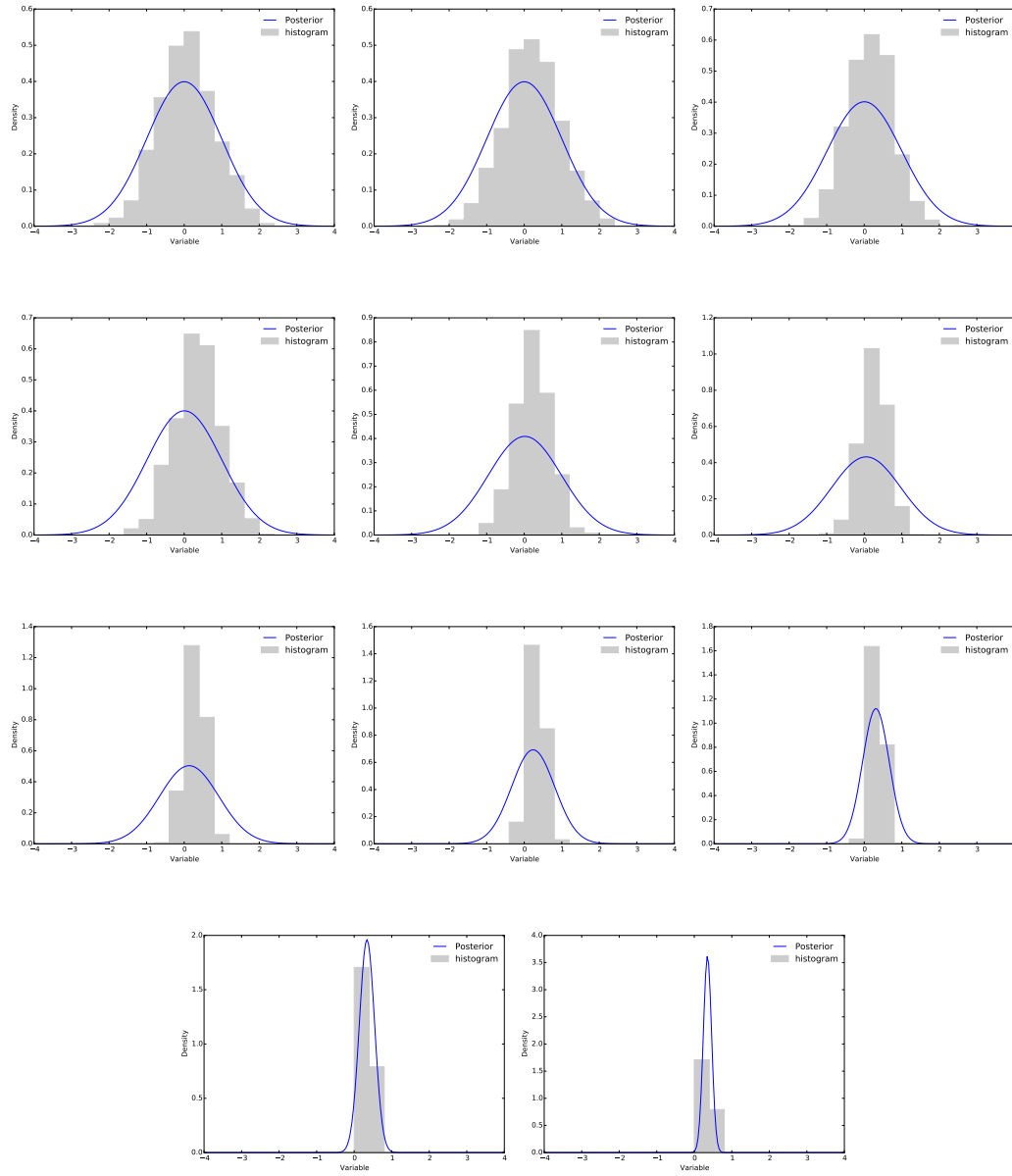
- The MCMC methods in Cotter et.al. [CRSW13] provide explicit examples of Markov kernels with the desired property of preserving measures, including the infinite dimensional setting.
- In fact, assuming stronger ergodicity properties for the move kernels, it is sometimes possible to obtain time-uniform bounds. See, e.g. the discussion at the end of previous chapter.

### 3.4 Adaptive SMC

In practice, the SMC samplers algorithm requires the specification of  $0 \leq \phi_0 < \phi_1 < \dots < \phi_p = 1$  as well as any parameters in the MCMC kernel. As demonstrated in Jasra et.al. [JSDT11], Kantas et.al. [KBJ13], the theoretical validity of which is established in Beskos et.al. [BJKT], these parameters can be set on the fly.

First, we focus on the specification of the sequence of distributions. Given step  $n - 1$  and  $\pi_{n-1}(x)$ , we select the next target density by adapting the temperatures to a required value of the effective sample size (ESS) (see Figure 3.1, Eq. 3.1) as in Jasra et.al. [JSDT11]. So, for a user defined threshold  $M_{thres}$ , we choose  $\phi_n$  as the solution of  $ESS_{(n)}(M) = M_{thres}$ . We use an inexpensive bisection search method to obtain the adapted temperature in this way.

To see how this adaptation behaves, we perform a simple experiment. Let the statistical model be a simple linear regression, i.e.  $Y = X + \varepsilon$  with  $\varepsilon \sim N(0, \sigma^2)$  and  $X \sim N(0, \theta^2)$ . We run an SMC with temperature adaptation as described above on this problem, with  $Y, X \in \mathbb{R}^{10}$ ,  $\sigma = 0.1$ ,  $\theta = 1$  and the number of particles is set to 1000. Our ESS target  $M_{thres}$  is 500. Figure 3.2 shows the resulting histograms of an arbitrarily selected component and the corresponding tempered targets. Note the smoothness of the transitions of targets. The distance between the consecutive densities should be small enough for the success of the algorithm.



**Figure 3.2:** Histograms of a single component of SMC particles at different iterations of the algorithm corresponding to the adaptively selected temperatures, together with the target tempered posteriors

Second, we turn to the specification of the move/mutation kernel  $K_n$ . Several options are available here, but we will use reflective random-walk Metropolis proposal on each component, conditionally independently. In particular, we will adapt the random move proposal scales  $\epsilon_{j,n}$  with  $j$  the coordinate and  $n$  the time index. A reasonable choice would be to adapt  $\epsilon_{j,n}$  to the marginal variance along the  $j$ -th coordinate; since this is analytically unavailable we opt for the SMC estimate from the previous time-step. Hence, using the notation from the previous chapter, we set  $\epsilon_{j,n} = \rho_n \sqrt{\hat{\text{Var}}_j(\xi_{n-1}^N)}$  where  $\rho$  is a global scaling parameter (recall that  $\xi_{n-1}^N$  denotes a random particle at time  $n-1$  within an  $N$ -particle system. We define  $\hat{\text{Var}}_j$  as the empirical variance operator of the  $j$ th component of the argument). For  $\rho$  itself, we propose to modify it based on the average acceptance rate at the previous time-step, aiming for an average acceptance rate around 0.2 (See Beskos et.al. [BRS09] for a theoretical justification). Our adaptive strategy halves (doubles)  $\rho$  if the last average acceptance rate went below (above) a predetermined neighbourhood of 0.2.

In addition, one can synthesize a number, say  $k_n$ , of baseline MCMC kernels to obtain an overall effective one with good mixing. To adapt  $k_n$ , we use the following heuristic: we propose to select  $k_n$  using  $k_n = \lfloor \frac{m}{\rho_n^2} \rfloor$  with  $m$  being a global parameter. The intuition is that for random walk type transitions of increment with small standard deviation  $\delta$ , one needs  $\mathcal{O}(\delta^{-2})$  steps to travel distance  $\mathcal{O}(1)$  in the state space. As a final modification, we enforce  $k_n$  to lie in a predetermined range, due to practical computational considerations.

The adaptive-SMC algorithm works as in Figure 3.1, except in step 1, before simulation from  $K_n$  is undertaken, our adaptive procedure is implemented. The algorithm will then run for a random number of time steps and terminate when  $\phi_n = 1$ , which will happen in finite time almost surely.

## 3.5 Numerical Results

### 3.5.1 Implementation details

The software used in our experiments has been implemented in C++ for the GNU/Linux platform. We used the Libmesh library for finite elements computation [KPSC06]. Fast

Fourier Transform was employed for rapid evaluation parameter  $u(x)$  at pre-determined grid-points and we exploited parallel computation wherever possible, via the MPI libraries. Our experiments were run on a computer server with 23 “Intel(R) Xeon(R)CPU X7460 @2.66GH” processors, each with 2 cores; 50 Gb memory and running “RedHat Linux version 2.6.18-194.el5” operating system. The experiments discussed in this paper used 20 processors. All the colour plots of random fields (e.g. permeability fields) have been prepared using the rainbow color scheme from the R programming language/environment. The scheme quantizes the Hue quantity of HSV (Hue Saturation Value) triplet of a pixel. Our level of quantization is selected to be 256 (8 bits), with the Hue range of  $[0, 1]$ , hence we normalize the random fields to this range and quantize to 8 bits to get the Hue value for a pixel. Saturation and Value were taken to be 1. All images were computed using  $500 \times 500$  equispaced point evaluations from the respective random fields.

### 3.5.2 Objects of inference

The work in [Vol13] investigates the performance of the Bayesian approach for our elliptic inverse problem and gives sufficient conditions under which posterior consistency holds. Posterior consistency is concerned with “ball probabilities” of type

$$\lim_{Card(O) \rightarrow \infty} \int_{B_\varepsilon} \frac{d\nu^y}{d\nu_0}(u) \nu_0(u) = 1$$

where  $y = \{y_x\}_{x \in O}$ ,  $O$  being the location set of the observations, and  $B_\varepsilon$  is the  $\varepsilon$  neighbourhood of the true value of  $u$ . One way to check such a result numerically is to use the posterior estimates obtained via our method. The estimated ball probabilities are computed as follows:

$$\sum_i w_p^i \mathcal{I}_{B_\varepsilon} u(x_p^i).$$

Although not all the conditions in [Vol13] required for posterior consistency to hold are fulfilled, we will nonetheless empirically investigate such a consistency property. This also provides a severe test for the SMC method since it implies posterior measures in the large dataset limit.

Parameter name	Value
frequency cutoff	10
finite elements d.o.f.	100
observation error std. dev.	$5 \times 10^{-7}$
number of Particles	1000
resampling threshold	600
wall-clock time	$\approx 11$ hrs

**Table 3.1:** Parameter values used for the 2D experiments. Between 5 and 1000 steps are allowed for the iterates of the MCMC kernels. The frequency cutoff determines the level of discretization of the permeability field. Finite elements d.o.f. denotes the number of finite elements used in the numerical solution of the elliptic PDE, higher values indicate better approximation at the expense of computational resources.

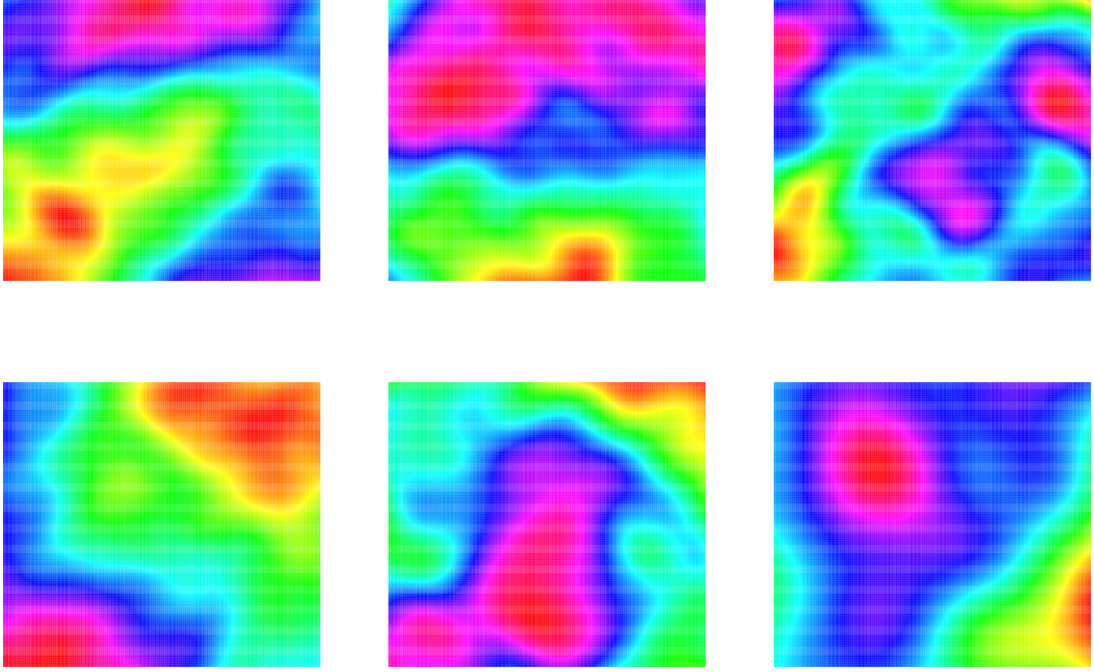
Before going into the results, we repeat our experimental setup for clarity. The observed random variable  $Y \in \mathbb{R}^d$  behaves according to  $Y = G(u) + \varepsilon$ , for  $\varepsilon \sim N(0, \sigma_\varepsilon I_{d \times d})$ . The operator  $G(u)$  is defined as  $G(u) = \mathcal{S}\mathcal{F}(u)$ , where  $\mathcal{F}$  is the forward solution of a PDE using FEM and  $\mathcal{S}$  is a sampling operator, sampling the underlying field on a regular grid. The goal is to compute the posterior density of  $u$  given observations  $Y$ .

### 3.5.3 2D Results

We consider the elliptic inverse problem in two dimensions. Our goal is to construct a sequence of posterior estimates, corresponding to increasing number of observations in order to numerically illustrate posterior consistency. Table 3.1 shows the parameters used in our experiments

To get an empirical sense of these parameters' effects on the distribution of the permeability field, we plot some samples from the prior field  $u(x)$  in Figure 3.3.

In another experiment, designed to study posterior consistency, a sequence of posterior estimates are formed by repeatedly running the adaptive SMC algorithm with, respectively, 4, 16, 36, 64 and 100 observations equi-spaced inside the 2D-domain. The computed MSE and ball probabilities are given in Figure 3.4, with the ball radius  $\varepsilon$  taken to be  $0.17 \times 380$ , where



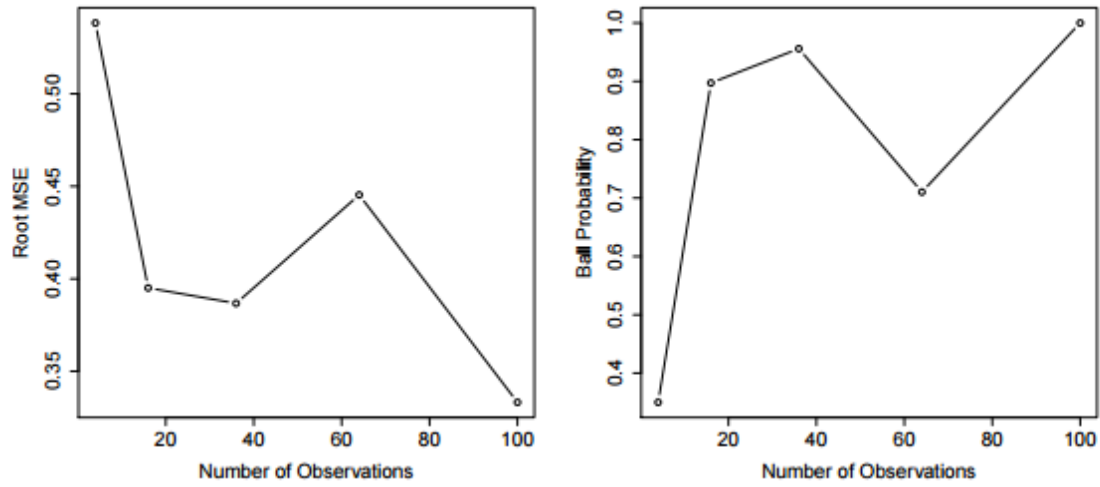
**Figure 3.3:** Six permeability field samples drawn from the prior

380 is the number of parameters in the system (i.e.  $2 \times 10 \times 19$ , 2 comes from the imaginary coefficients of the Fourier transform, 10 is the cardinality of  $\{0, \dots, 9\}$  and 19 corresponds to  $\{-9, \dots, 0, \dots, 9\}$ ), corresponding to a frequency cutoff of 10. The Figure suggests that as more data become available posterior consistency is obtained as predicted, under slightly more restrictive assumptions than we have in play here, in [Vol13]. This is interesting for two reasons: firstly it suggests the potential for more refined Bayesian posterior consistency analyses for nonlinear PDE inverse problems; secondly it demonstrates the potential to solve hard practical Bayesian inverse problems and to obtain informed inference from a relatively small number of observations.

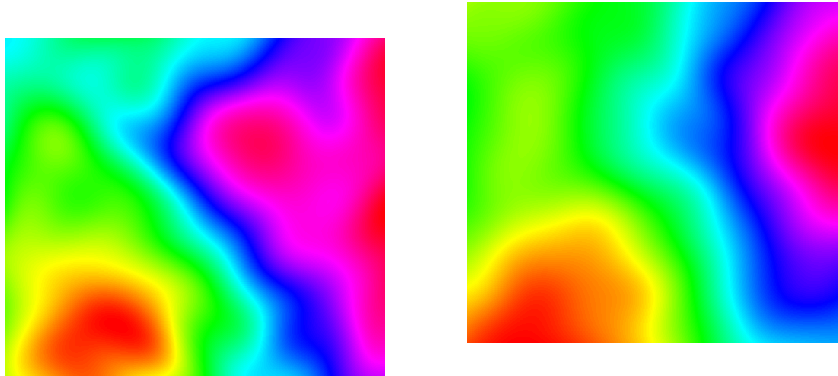
Figure 3.7 shows an example inference, with the estimated permeability field on the right, and the corresponding true field on the left.

Finally, Figure 3.8 shows marginal posterior density estimates corresponding to 144 observations. The usual observation is to note the effectiveness of even the mode estimator in





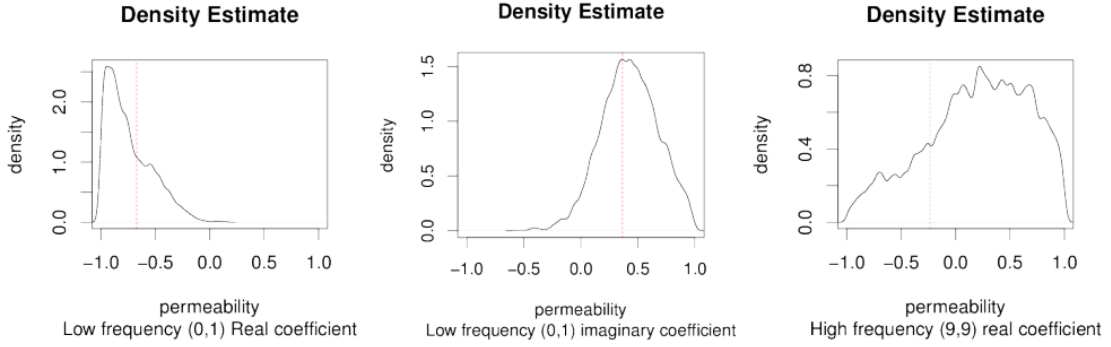
**Figure 3.4:** Numerical consistency checks for the sequence of experiments with 4,16,36,64 and 100 observations



**Figure 3.5:** True Permeability Field  
**Figure 3.6:** Estimated Permeability

**Figure 3.7:** An estimated permeability field and the corresponding true field

lower frequencies. Another important observation is the similarity of the high frequency marginal densities to the prior. In fact, it is this behaviour that makes a prior invariant MCMC proposal superior to others, i.e. the proposal itself is almost optimal for a wide range of coefficients in the problem.



**Figure 3.8:** Posterior marginal density estimates for two low and one high frequency coefficients in the 2D case

Parameter name	Value
number of observations	125
frequency cutoff	5
finite elements d.o.f.	1000
observation error std. dev.	$1 \times 10^{-8}$
number of Particles	1000
resampling threshold	600
wall-clock time	$\approx 10$ days

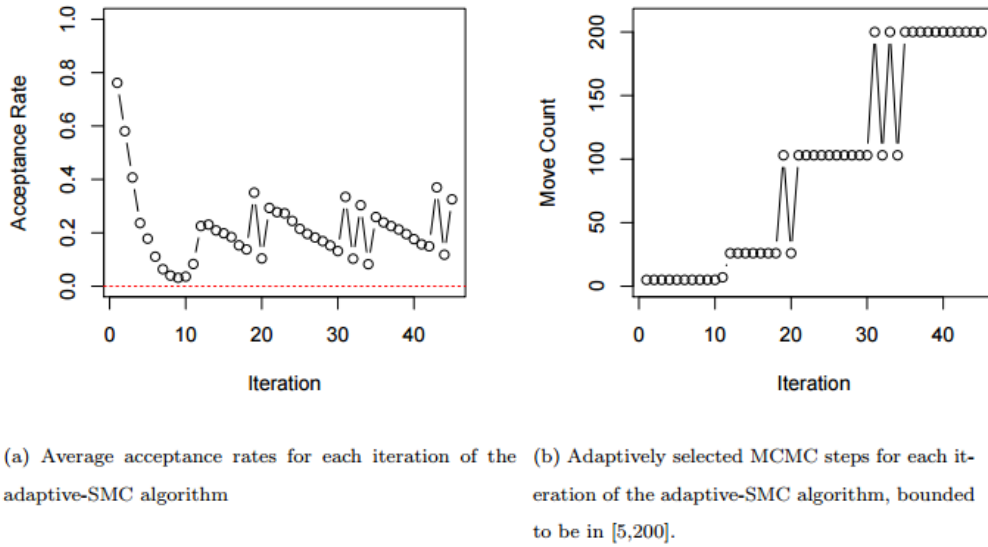
**Table 3.2:** Parameter values used for the 3D experiment. Between 5 and 200 steps are allowed for the iterates of the MCMC kernels.

### 3.5.4 3D Results

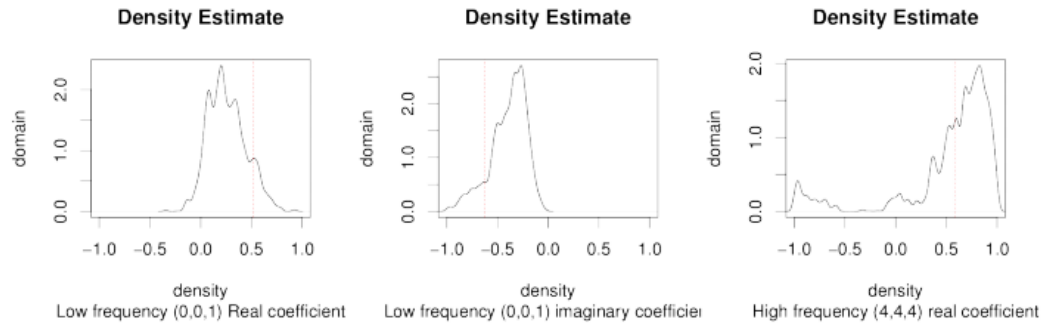
A more realistic experiment is performed using the 3-D setup. In this setup, the computational aspects of the problem are further highlighted as the numerical solution of the forward operator becomes much harder due to the increased cardinality of the finite elements basis. The values of parameters in this numerical study are given in Table 3.2. The data are generated from the model, under the specifications given in Table 3.2.

In Figure 3.9, we consider the performance of our SMC algorithm in this very challenging scenario. In Figure 3.9 (a), we can see the average acceptance rates of the MCMC moves

over the time parameter of the SMC algorithm. We can observe that these acceptance rates do not collapse to zero and are not too far from 0.2. This indicates that the step-sizes are chosen quite reasonably by the adaptive SMC algorithm and the MCMC kernels have some mixing ability. In Figure 3.9 (b), we can see the number of MCMC iterations that are used per-particle over the time parameter of the SMC algorithm. We can observe, as one might expect, that as the target distribution becomes more challenging, the number of MCMC steps required grows. Figure 3.9 indicates reasonable performance of our SMC algorithm. In terms of inference, the posterior density estimates are shown in Figure 3.10. Recall that the priors are uniform. These estimates indicate a clear deviation from the prior specification, illustrating that the data influence our inference significantly. This is not obvious, and establishes that one can hope to use this Bayesian model in real applications.



**Figure 3.9:** SMC Performance for 3D Example.



**Figure 3.10:** Posterior marginal density estimates for two low and one high frequency coefficients in the 3D case

# CHAPTER 4

## Multi-Resolution Sequential Monte Carlo Inference

In the last chapter, we attempted to solve the inverse problem by considering it as a generic inference problem with computationally complex likelihood evaluation. We saw that, even with novel ways of allocating the computational resources, the resulting algorithm is prohibitively expensive, where a relatively small size problem can take days of run time even with a fully parallel implementation. Hence, we now consider the particular structure of the underlying problem, using multi-resolution Monte Carlo inference ideas in a SMC setting. This chapter is based upon our project of investigating the practicality of SMC in applied inverse problems. The work is based on [BJL<sup>+</sup>15], where the authors perform the relevant error-analysis and showcase the resulting algorithm in a toy problem. Our work is a continuation of theirs, in the sense that we investigate the practical properties of this algorithm in a realistic groundwater-flow problem. This is the second chapter in this thesis that incorporates our novel contributions. We begin the chapter with a discussion of the related idea of using fast approximations of the underlying posterior to speed up inference.

### 4.1 Speeding up Monte Carlo Computations in Inverse Problems

Bayesian approach to inverse problems involves at least three approximations: discretization of the parameter field, Galerkin projection of the forward operator and the Monte Carlo approximation to the expectations. Hence, from the perspective of inverse problems, even a perfect evaluation of the expectation of interest would yield an approximate solution, since the underlying posterior is only an approximation. On the other hand, the error components

coming from the approximation of the posterior are out of the reach of the Monte Carlo method. Hence, from the perspective of Monte Carlo inference, the methods discussed so far are exact (in the sense that they can evaluate the expectation with arbitrarily small error). Some of the methods we will discuss in this section introduce additional approximations, in order to cut computational costs. These approximations may be corrected eventually or they may end up adding another error component to the result. All of the methods discussed here will attempt to reduce the high computational cost of “plain” Monte Carlo methods, as exemplified in the previous chapter. We will see that significant gains are possible when the underlying model structure is exploited. This discussion will also form a bridge to the idea of multi-resolution Monte Carlo methods, which enable even further gains.

#### 4.1.1 Early Rejection

We begin our discussion with the work of Solonen et.al. [SOL<sup>+</sup>12]. They present two distinct approaches to speed up inference, which can be implemented together. Their goal is to make inference in complicated climate models [RfM03], with high-cost likelihood evaluations. They report that normally such models allow efficient parallelization for only about a few dozen number of processors. Given the high cost of inference in such models, additional (efficient) parallelization is desired. To achieve this, they consider using a batch of parallel MCMC chains that interact [CRY09]. Furthermore, adaptive scaling of the Metropolis-Hastings proposal is adopted to increase the acceptance rate [HST01]. The resulting adaptive parallel MCMC method uses a symmetric random walk proposal whose covariance matrix is adapted at regular intervals using the output of all chains. This asynchronous update mechanism works as follows; let  $U$  be the parameter of interest,  $\bar{U}_i$  (adapted mean at iteration  $i$ ) and  $\Sigma_i$  (adapted covariance at iteration  $i$ ) initially set to their prior values. We define  $i$  as a global counter, initially set to 0. When the chain  $j$  finishes its iteration, we perform the update  $\bar{U}_{i+1} = \bar{U}_i + \frac{1}{1+i}(U^j - \bar{U}_i)$ , where  $U^j$  is the state of the chain  $j$ . We then update the covariance as  $\Sigma_{i+1} = \frac{i-1}{i}\Sigma_i + \frac{1}{i}(U^j - \bar{U}_i)(U^j - \bar{U}_i)^T$  and increase the counter  $i = i + 1$ . We observe that the adaptations become less impactful as the time goes, which is a common feature in adaptive MCMC.

Application of MCMC methods involve a substantial amount of simulation, which are very costly in inverse problems. The second idea proposed in [SOL<sup>+</sup>12] is to stop the computations early when the rejection is certain. They achieve this by reversing the order of simulation and likelihood evaluation in a typical MCMC step; where normally one evaluates the likelihood, leading to the Metropolis-Hastings ratio which is compared to a value sampled from the uniform distribution. Assume the posterior is factorized as  $\pi(u) \propto \pi_0(u) \prod_{i=1}^n L(u; y_i)$ , where  $L(u; y_i)$  is the likelihood for the  $i$ th observation  $y_i$ , which are assumed to be conditionally independent. Define  $\pi_k(u) = \pi_0(u) \prod_{i=1}^k L(u; y_i)$ , and observe that it is monotonically decreasing if the likelihood is bounded (which is the case when  $L(u; y_i) \propto \exp(-\Phi(u; y_i))$ ) so that it can be normalized to be  $|L| < 1$ <sup>1</sup>. Since we use a symmetric random walk proposal, the Metropolis Hastings ratio is  $\frac{\pi(u_{t+1})}{\pi(u_t)} > \frac{\pi_k(u_{t+1})}{\pi_k(u_t)}$ . The algorithm proceeds by generating a  $v \sim U[0, 1]$  and comparing  $\frac{\pi_k(u_{t+1})}{\pi_k(u_t)}$  to  $v$  sequentially, rejecting as soon as the MH ratio falls below  $v$ . Considerable savings are possible when the cost of per-observation-likelihood evaluation is significant. However, in our problem, the real cost is the forward solution and the consecutive computations for per-observation-likelihoods are negligible. It turns out, however, that the related idea of delayed acceptance is much more relevant.

#### 4.1.2 Delayed Acceptance

A related idea to early rejection is that of delayed acceptance, originally proposed in [FN97, CF05]. The goal is to use a surrogate (approximate) density  $\hat{\pi}$  to perform an initial accept-reject step, using the acceptance probability  $\hat{\rho}(x, y) = \min \left\{ \frac{q(y, x) \hat{\pi}(y)}{q(x, y) \hat{\pi}(x)}, 1 \right\}$ . If the new state is accepted in this initial stage, another accept-reject step is performed, this time targeting the actual density, i.e.  $\rho(x, y) = \min \left\{ \frac{q(y, x) \hat{\rho}(y, x) \pi(y)}{q(x, y) \hat{\rho}(x, y) \pi(x)}, 1 \right\}$ . This two stage process has equal or smaller acceptance rate for all state transitions, as we will show when we discuss the general method, and hence has worse asymptotic variance as per Peskun's ordering of Markov chains [Pes73]. Yet, in certain problems and with a good choice of surrogate density this method can mix more efficiently with respect to the computational resources.

A generalization of this idea as well as early rejection is proposed by Banterle et.al. [BGLR15].

---

<sup>1</sup>recall that we are not concerned with normalization constants, as they are cancelled in the MH ratio

The idea is to factorize the Metropolis-Hastings ratio as  $r(x, y) = \Pi_k r_k(x, y)$  such that  $r_k(x, y) = r_k^{-1}(y, x)$ , where the factors may correspond to prior, a factor of likelihood or some surrogate density. The key observation is that  $\widehat{\rho}(x, y) = \Pi_k \min\{r_k(x, y), 1\}$  can be used as an acceptance probability to create a chain that targets  $\pi$ , i.e.  $\widehat{P}(x, A) = \int_A \widehat{\rho}(x, y) Q(x, dy) + \left( \int_X (1 - \widehat{\rho}(x, y)) Q(x, dy) \right) \mathbf{1}_A(x)$ . To see this, we re-express the detailed balance condition as  $\frac{\rho(x, y)}{\rho(y, x)} = r(x, y) = \frac{q(y, x)\pi(y)}{q(x, y)\pi(x)}$ . Therefore, we need to show that  $\frac{\widehat{\rho}(x, y)}{\widehat{\rho}(y, x)} = r(x, y)$ .

$$\begin{aligned} \frac{\widehat{\rho}(x, y)}{\widehat{\rho}(y, x)} &= \frac{\Pi_k \min\{r_k(x, y), 1\}}{\Pi_k \min\{r_k(y, x), 1\}} \\ &= \Pi_k \frac{\min\{r_k(x, y), 1\}}{\min\{r_k(y, x), 1\}} \\ &= \Pi_k r_k(x, y) = r(x, y), \end{aligned} \tag{4.1}$$

since  $r_k(x, y) = r_k^{-1}(y, x)$  and  $\frac{\min(a, 1)}{\min(a^{-1}, 1)} = a$  for any positive  $a$ . We also observe that,

$$\widehat{\rho}(x, y) = \Pi_k \min\{r_k(x, y), 1\} \leq \min\{\Pi_k r_k(x, y), 1\} = \min\{r(x, y), 1\} = \rho(x, y),$$

since  $\min(a, 1) \min(b, 1) \leq \min(ab, 1)$  for  $a, b \in \mathbb{R}^+$ . As discussed before, Peskun's ordering implies that asymptotic variance of the delayed acceptance Markov chain (with kernel  $\widehat{P}$ ) is larger than that of the standard Metropolis-Hastings method (with kernel  $P$ ). Furthermore, we have the following result,

**Theorem 18.** [BGLR15] Define  $A = \{(x, y) | r(x, y) > 1\}$ . Assume there exists  $c > 0$  such that,

$$\inf_{(x, y) \in A} \min_k r_k(x, y) \geq c. \tag{4.2}$$

Then we have,

$$\text{Var}(f, \widehat{P}) \leq (c^{1-d} - 1) \text{Var}_\pi(f) + c^{1-d} \text{Var}(f, P) \quad f \in L_0^2$$

with  $L_0^2 = \{f | \pi(f) = 0 \text{ and } \pi(f^2) < \infty\}$ , and

$$\text{Var}(f, Q) = \lim_{n \rightarrow \infty} \text{Var} \left( n^{1/2} \sum_{i=1}^n (f(X_i) - \pi(f)) \right)$$



where  $(X_i)$  is a Markov chain that is initialized as  $X_1 \sim \pi$ .

An important corollary is that, if the standard Metropolis Hastings MCMC is geometrically ergodic, then, assuming Eq. 4.2, the delayed acceptance MCMC retains this property. Banterle et.al. [BGLR15] also shows that if this assumption does not hold, it is possible to construct examples where the standard Metropolis-Hastings is geometrically ergodic but the delayed acceptance algorithm is not. A slight modification overcomes this issue, if we set

$$\hat{r}_k(x, y) = \min \{1/b, r_k(x, y)\}$$

and,

$$\hat{r}_d(x, y) = \frac{r(x, y)}{\prod_{k=1}^{d-1} \hat{r}_k(x, y)},$$

so that we still have  $r(x, y) = \prod_{k=1}^d \hat{r}_k(x, y)$ . This setup guarantees that the assumption of Eq. 4.2 holds, therefore the resulting chain is geometrically ergodic.

We end this discussion with an example of this procedure [CF05] that employs surrogate densities, which is very useful in inverse problems. Suppose that we have a surrogate posterior  $\hat{\pi}$  that approximates the posterior  $\pi$  but is less costly to compute. In inverse problems, this can correspond to a posterior obtained using a small resolution FEM solution in likelihood evaluations. We assume a symmetric proposal, so that the Metropolis-Hastings ratio reads  $r(x, y) = \frac{\pi(y)}{\pi(x)}$ . Let the first factor be  $r_1(x, y) = \frac{\hat{\pi}(y)}{\hat{\pi}(x)}$ , and the second factor  $r_2(x, y) = \frac{\hat{\pi}(x)\pi(y)}{\hat{\pi}(y)\pi(x)}$ . We observe that this is the same procedure as the one discussed in the beginning of this subsection. In the implementation one would first perform an accept-reject step with respect to the acceptance probability of  $\min\{r_1, 1\}$  which is presumably very easy to evaluate and then perform the second test (with acceptance probability  $\min\{r_2, 1\}$ ) only if the first one is accepted. Significant speedup can be obtained if the first test correlates highly with the standard Metropolis-Hastings test.

### 4.1.3 Surrogate Models

Up to now, we have only considered approximate distributions based on various resolutions of the numerical solution of the forward problem. These are far from being the only choices,

when one wants to replace a complicated density with an approximate one that is easier to compute. In fact, a whole literature exists around constructing such densities, named “surrogate models”, “meta-models” or “statistical emulators” [KS15, HSS12, MN09, CMW15]. Such approximate densities can be used as part of, for example, delayed acceptance MCMC or other such methods that eventually correct these approximations, but most often they are simply used as new approximate targets. We will now briefly discuss polynomial chaos expansions [MN09], low-rank polynomial approximations [KS15] and a data driven collocation method [CMW15]. The common property of all these methods is that they involve a costly initial parameter determination phase, but once the parameters of the approximations are determined, approximate forward solutions become very cheap to compute. Hence, these methods are suited for applications where one expects to compute a lot of forward solutions, such as long MCMC runs.

Polynomial chaos expansion relies on the fact that any square-integrable random variable  $X$  can be expressed as [MN09],

$$X = a_0\Gamma_0 + \sum_{i_1} a_{i_1}\Gamma_1(\xi_{i_1}) + \sum_{i_1, i_2} a_{i_1 i_2}\Gamma_2(\xi_{i_1}, \xi_{i_2}) + \dots$$

where  $\Gamma_p$  is the polynomial chaos of order  $p$ . This expression can be rewritten as,

$$X = \sum_{i=1}^{\infty} a'_i \Psi_i(\xi_1, \xi_2, \dots)$$

such that each  $a'_i$  corresponds to an  $a_\alpha$  as well as each basis  $\Psi_i(\xi_1, \xi_2, \dots)$  to one of  $\Gamma_p$ . We require that  $\Gamma_p$  (and hence  $\Psi_i$ ) form an orthogonal basis and that  $\xi_i$  are i.i.d. standard normal variables. The exact functional form of  $\Gamma_p$  depends on the distribution of  $X$ , for example a uniform variable on  $[-1, 1]$  corresponds to the family of Legendre polynomials, whereas a standard normal  $X$  is associated with Hermite polynomials [KS15]. Such an expansion can, of course, only be practically utilized after truncation of both the number of variables  $\{\xi_i\}_i$  and the order of the polynomials. If we choose the truncation such that  $\text{Card}\{\xi_i\}_i = n$  and the highest order of the polynomials is  $p$ , we obtain an expansion with

$P = \binom{n+p}{n}$  [MN09],

$$X \approx \sum_{i=1}^P a'_i \Psi_i(\xi_1, \xi_2, \dots, \xi_n).$$

The goal in inverse problems is to use polynomial chaos expansion on  $X = G(U)$ , where  $G(\cdot)$  is the forward solution and  $U$  is the unknown field. The idea is to replace FEM computations for each forward solution with a simple polynomial approximation, which is much faster to evaluate. The coefficients  $a'_i$  of this approximation can be determined in different ways. The *intrusive method* [MN09] directly employs the underlying model leading to the forward operator  $G$ . In the general case, if we have an operator  $\mathcal{O}(X, U) = 0$ , we simply replace  $X$  and  $U$  with their PC-expansions and solve the resulting linear system for the PCE coefficients. For example, in the case of ground-water problem, this operator is  $-\nabla \cdot (\kappa \nabla U) - f = 0$  ;  $x \in D$ , under the boundary condition  $U = 0$  ;  $x \in \partial D$ . We note the similarity of this approach to Galerkin projection described in Chapter 1. On the other hand, the *non-intrusive method* considers the underlying problem as a black box. The problem is modeled like a standard regression problem, i.e.  $X = \sum_{i=1}^P a'_i \Psi_i(U) + \varepsilon$ , where  $\varepsilon$  corresponds to the truncated terms. A least squares solution can be obtained by minimizing an empirical mean squared error,

$$\mathbf{a}' = \arg \min_{\mathbf{c}} \sum_j ||G(u_j) - \sum_{i=1}^P c_i \Psi_i(u_j)||^2,$$

where  $u_i$  may either be predetermined values, or sampled from standard normal. The solution can be expressed as  $\mathbf{a}' = (\Psi^T \Psi)^{-1} \Psi^T \mathcal{X}$ , where  $\Psi_{ij} = \Psi_j(u_i)$  and  $\mathcal{X}_i = G(u_i)$ . As we discussed before, the majority of the computational budget goes towards computing  $\mathcal{X}$  with a FEM solution for each component.

A similar polynomial expansion, under the name of low rank approximations, is proposed in [KS15]. The proposed approximation is of the form,

$$X \approx \sum_{l=1}^R b_l (\Pi_{i=1}^M v_l(X_i)),$$

where  $v_l(\cdot)$  are univariate functions constructed using univariate polynomials  $P_k$  up to order

$p_i$ , i.e.  $v_l(X_i) = \sum_{k=0}^{p_i} z_{k,l} P_k(X_i)$ . The parameters of this approximation are  $b_l$ 's and  $z_{k,l}$ 's. Therefore the number of parameters are  $R \cdot (\sum_{i=1}^n (p_i + 1))$ . We observe that this number grows linear with respect to  $n$  (the input dimensionality) as opposed to the exponential growth of PCE coefficients. This forms one of the motivations of this approach, assuming good approximations can be constructed this way. The parameters of this approach are determined in a greedy fashion in [KS15], where for each parameter, the others are assumed fixed and the resulting single variable least squares problem is solved. Several heuristics to find an appropriate  $R$ , called the rank of the approximation, is proposed in [KS15]. This is another example of a non-intrusive method.

As we have seen so far, the non-intrusive approach computes the forward map at certain points  $\{u_i\}_i$ , which are either sampled from standard normal or selected deterministically, and then finds the expansion coefficients that most closely approximate these points in the MSE sense. The hope is that, the resulting polynomial approximation will also be a reasonably close approximation for other points as well. An important observation is that, the forward map  $G(\cdot)$  at points  $u$  that are close to the “snapshot” set  $\{u_i\}_i$  is better approximated than at other points. Based on this, Cui et.al. [CMW15] propose to sample  $\{u_i\}_i$  from posterior, leading to better approximations of  $G(\cdot)$  in parts of the space that matters most. They also show that, in this case, using simple linear basis (i.e. selecting the maximum order  $p$  as 1) functions suffices. They use their algorithm within a delayed acceptance scheme and show that significant computational savings can be obtained.

## 4.2 Multi-resolution Monte Carlo

Inference in inverse PDE problems often involve a discretization of the underlying continuum field, such as finite elements method using Galerkin projection. Various levels/resolutions of refinement of such discretizations correspond to a tradeoff between accuracy and computational complexity; giving a natural hierarchy of resolutions. A practitioner selects a discretization level by matching a desired error level; however computations at very fine resolutions may often be prohibitively complex. A big literature exists for the solution to this common problem: how can we utilize the computations at coarser resolutions (e.g.

as surrogate models) to be able to achieve a given error with much less computational resources? For example, to solve a linear system, one can use the solution corresponding to a lower resolution as a preconditioner for a higher resolution solvers. This is the principle of multi-grid methods. In the context of Monte Carlo methods, multi-resolution methods can be viewed in two ways: for a fixed computational cost these methods are variance reduction methods; on the other hand, for a fixed target error, these methods reduce the computational complexity. This effect is achieved via a telescoping sum associated with the hierarchy of resolutions.

We now make this idea more concrete under an ideal i.i.d. Monte Carlo sampling scheme. The target quantity is an expectation of a functional  $g$  of the parameter of interest  $U$ . We denote the ideal (non-discretized) measure of  $U$  as  $\nu_\infty$ . The approximate law corresponding to resolution  $l$  is denoted as  $\nu_l$  and the corresponding density  $\pi_l$ . This method uses the following telescoping sum

$$E_{\pi_L}(g(u)) = E_{\pi_0}(g(u)) + \sum_{i=1}^L (E_{\pi_i}(g(u)) - E_{\pi_{i-1}}(g(u))).$$

The main observation is that,  $[E_{\pi_i}(g(u)) - E_{\pi_{i-1}}(g(u))]$  diminishes as  $i \rightarrow \infty$ . This hints at the possibility that, at higher resolutions, one can hope to use less computational resources and still get a reasonable estimator. A Monte Carlo estimator for each term of the telescoping sum is as follows

$$Y_l^{N_l} = \sum_{i=1}^{N_l} (g(u_l^{(i)}) - g(u_{l-1}^{(i)})) N_l^{-1}$$

with  $g(u_{-1}) := 0$ . Here  $\{u_l^{(i)}, u_{l-1}^{(i)}\}$  are i.i.d. with marginal densities  $\pi_l, \pi_{l-1}$ . Hence the overall estimator reads

$$Y_{L,Multi} = \sum_{l=0}^L Y_l^{N_l}.$$

A simple decomposition of the mean squared error gives;

$$E(Y_{L,Multi} - E_{\pi_\infty}(g))^2 = E(Y_{L,Multi} - E_{\pi_L}(g))^2 + (E_{\pi_L}(g) - E_{\pi_\infty}(g))^2 \quad (4.3)$$

where the first term is variance and the second is bias-squared. The cross term vanishes

because we have an unbiased estimator. Furthermore, since the samples  $\{u_l^{(i)}, u_{l-1}^{(i)}\}$  are i.i.d., the variance term decomposes into

$$V_l = \text{Var}(g(u_l^{(i)}) - g(u_{l-1}^{(i)})).$$

The key observation here is that  $V_l$  decreases rapidly as the differences converge to 0. This implies that, if we target a certain error for the variance component of the MSE, we may need much less particles at higher levels. In fact, we can calculate exactly how the  $N_l$  should decay. For this, take  $C_l$  as the computational cost of one sample at resolution  $l$ . One can calculate the optimal allocation of particle counts  $((N_l)_{l=0}^N)$  for each resolution by minimizing  $\sum_l V_l/N_l$  for fixed  $\sum_l C_l N_l$ . Lagrange multipliers method yields  $N_l \propto \sqrt{V_l/C_l}$ . The only other remaining parameter, which is the number of resolutions  $L$ , can be calculated by matching the bias to the desired order.

### 4.3 Multi-resolution Monte Carlo Path Simulation

In this section we will discuss the famous Multi-resolution Monte Carlo path simulation example by Giles [gil08]. This example successfully puts our previous discussion into perspective, as well as showing a concrete performance gain. Path simulation problem is concerned with the estimation of quantities like  $E(g(S(T)))$ , where  $S(t)$  is a random process (and  $S(T)$  its terminal state) defined by a Stochastic Differential Equation (SDE). They use the following model with drift and volatility terms,

$$dS(t) = a(S, t)dt + b(S, t)dW(t), \quad 0 < t < T,$$

and some given initial data  $S_0$ . In the above,  $W$  denotes the Wiener process. In applications, such a process is discretized first, which enables the simulation of paths (and hence the terminal states). An Euler discretization of this SDE with timestep  $h$  is

$$\hat{S}_{n+1} = \hat{S}_n + a(\hat{S}_n, t_n)h + b(\hat{S}_n, t_n)\Delta W_n$$

$$t_{n+1} = t_n + h$$

where  $\Delta W_n$  is a Gaussian random variable, corresponding to the  $h$ -jump of the Wiener process at  $t_n$ . This discrete approximation of the process enables easy path sampling of  $\widehat{S}_n^{(i)}$ , which is then used in the following estimator:

$$\widehat{Y} = N^{-1} \sum_i g(\widehat{S}_{T/h}^{(i)}).$$

Provided that the drift  $(a(S, t))$  and the volatility  $(b(S, t))$  functions satisfy certain conditions [BT, gil08], the expected mean squared error (MSE) of the estimate  $\widehat{Y}$  is asymptotically of the form  $MSE \approx c_1 N^{-1} + c_2 h^2$ , similar to the discussion of the last section. An MSE of order  $O(\varepsilon^2)$  is obtained by setting  $N = O(\varepsilon^{-2})$  and  $h = O(\varepsilon)$ , resulting in a computational complexity of  $C = Nh^{-1} = \varepsilon^{-3}$ . The multi-resolution method will reduce this complexity to  $O(\varepsilon^{-2}(\log(\varepsilon))^2)$ , using coarser path simulations to reduce variance, while employing a correction that retains the bias of the finest resolution.

We begin by defining a refinement strategy, i.e. choosing  $h_l$  as  $M^{-l}T$  for  $l = 0, 1, \dots, L$ , where  $M$  is a positive real number. Let  $P = g(S(T))$  and let  $\widehat{P}_l$  denote its approximation at discretization level  $l$ . Finally we denote by  $\widehat{P}_l^{(i)}$  the  $i$ th sample of  $\widehat{P}_l$ .

We begin the development of the multi-resolution estimator by re-expressing  $E(\widehat{P}_L)$  as:

$$E(\widehat{P}_L) = E(\widehat{P}_0) + \sum_{i=1}^L E(\widehat{P}_i - E(\widehat{P}_{i-1})).$$

The strategy is to estimate each term independently, in a way that will minimize the computational complexity for a target error. Let  $\widehat{Y}_l = N_l^{-1} \sum_i \widehat{P}_l^{(i)}$  be the Monte Carlo estimators for  $\widehat{P}_l$ . We estimate the first term in the above telescoping sum by  $\widehat{E}(\widehat{P}_0) = \widehat{Y}_0$  and the rest by

$$\widehat{Y}_l = N_l^{-1} \sum_i (\widehat{P}_l^{(i)} - \widehat{P}_{l-1}^{(i)}).$$

The key point here is to use the same sample path  $\widehat{S}_{n,l}$  for both  $\widehat{P}_l^{(i)}$  and  $\widehat{P}_{l-1}^{(i)}$ . This can easily be achieved by sampling the finer path first, and then summing the corresponding  $\Delta W_n$

realizations in  $M$  groups. For example, if  $S_2$  is calculated using  $\Delta W_{2,n} = [0.2, 0.3, -1, 1.2]$  then  $S_1$  can be calculated using  $\Delta W_{1,n} = [0.2 + 0.3, -1 + 1.2]$ , assuming  $M$  is 2. This means only  $N_l$  samples are generated for each  $\hat{Y}_l$ . Since the paths are independent, the variance of this simple estimator is  $V(\hat{Y}_l) = N_l^{-1}V_l$ , where  $V_l$  is the variance of a single sample. Notice that the  $\hat{Y}_l$  are also independent for different  $l$ 's, therefore the variance of the overall estimator is

$$V(\hat{Y}) = \sum_l N_l^{-1}V_l.$$

On the other hand, the computational cost is proportional to  $\sum_l N_l h_l^{-1}$ . As we have seen in the last section, Lagrange multipliers method yields  $\sqrt{V_l h_l}$  as the optimal number of samples per level. Now we consider the case of  $L \gg 1$  and analyze the behaviour of  $V_l$  as  $l \rightarrow \infty$  and, in turn, that of  $V$ . For the case of Euler discretization, under certain assumptions [gil08, BT], we have the following strong convergence result

$$E(\|\hat{S}_l - S(T)\|^2) = O(h_l).$$

Assuming the function  $g(\cdot)$  is Lipschitz, we have

$$V(\hat{P}_l - P) \leq CE(\|\hat{S}_l - S(T)\|^2).$$

Combining these two results, we have  $V(\hat{P}_l - P) = O(h_l)$ . Additionally, Minkowski's inequality gives:

$$V(\hat{P}_l - \hat{P}_{l-1}) \leq \left( (V(\hat{P}_l - P))^{1/2} - (V(\hat{P}_{l-1} - P))^{1/2} \right)^2.$$

Hence we find  $V_l = O(h_l)$ , and the optimal choice of  $N_l$  is asymptotically  $h_l$ . Putting this into the expression  $V(\hat{Y}) = \sum_l N_l^{-1}V_l$ , we find that setting  $N_l = O(\varepsilon^{-2}Lh_l)$  results in  $V = O(\varepsilon^2)$ . Now it remains to find the  $L$  that give a bias with the same order. If we choose  $L = \frac{\log \varepsilon^{-1}}{\log M} + O(1)$  as  $\varepsilon \rightarrow 0$ , then  $h_L = M^{-L} = O(\varepsilon)$ , and we get the required bias. If we combine  $N_l = O(\varepsilon^{-2}Lh_l)$  with the expression for computational cost  $\sum_l N_l h_l^{-1}$ , we get the



asymptotic cost of  $O(\varepsilon^{-2}(\log \varepsilon)^2)$ .

## 4.4 Multi-resolution Sequential Monte Carlo

In our elliptic inverse problem setup, a difficulty arises since we can only evaluate up to a constant the target density and cannot directly obtain independent samples from it. The most common approach to solve this problem in the literature is to apply MCMC methodology [HSS12]. On the other hand, a recent article [BJL<sup>+</sup>15] showed how to apply SMC to this end and developed relevant theory. In this section, we will overview their developments. Later this chapter we will present empirical results based on our own implementation.

The SMC methodology applied to this problem works essentially the same as in the previous chapter. The difference is that the sequence of target distributions are now posteriors corresponding to an increasing sequence of mesh resolutions. Contrasting this to the i.i.d. Monte Carlo approach outlined previously, we see that  $\{u_l^{(i)}, u_{l-1}^{(i)}\}$  are no longer i.i.d. but are dependent particles evolved from the prior using SMC. This also implies that the  $Y_l^{N_l}$  are no longer unbiased estimators of  $E_{\pi_l}(g) - E_{\pi_{l-1}}(g)$ . In turn, this means that the previous decomposition of MSE in Eq. 4.3 is no longer relevant. Instead we have,

$$E((Y - E_{\pi_\infty}(g))^2) \leq 2E((Y - E_{\pi_L}(g))^2) + 2(E_{\pi_L}(g) - E_{\pi_\infty}(g))^2.$$

To develop multi-resolution SMC method, we begin by re-writing the telescoping sum as,

$$E_{\pi_0}(g(u)) + \sum_l E_{\pi_{l-1}} \left( \left( \frac{\pi_l(u)}{\pi_{l-1}(u)} - 1 \right) g(u) \right).$$

We evolve a particle system targeting  $\pi_0, \pi_1, \dots, \pi_L$  using SMC sampler, i.e.  $\{u_0^{1:N_0}, u_1^{1:N_1}, \dots, u_{L-1}^{1:N_{L-1}}\}$ .

We resample  $u_l^{1:N_l}$  with weights  $G_l(u_l) = (\pi_{l+1}/\pi_l)(u_l)$ . Finally, the resampled particles are moved with kernels  $K_l$  s.t.  $\pi_l K_l = \pi_{l+1}$ . This particle system is the basis of the SMC estimator of  $E_{\pi_{l-1}}(\phi(u))$  :

$$\pi_{l-1}^{N_{l-1}}(\phi) = \frac{1}{N_{l-1}} \sum_{i=1}^{N_{l-1}} \phi(u_{l-1}^{(i)})$$

The overall ML-SMC estimator becomes;

$$Y = \pi_0^{N_0}(g) + \sum_l \left( \frac{\pi_{l-1}^{N_{l-1}}(gG_{l-1})}{\pi_{l-1}^{N_{l-1}}(G_{l-1})} - \pi_{l-1}^{N_{l-1}}(g) \right).$$

The dependency structure of SMC also means that the optimal number of particles is harder to calculate. In particular,  $Y_l^{N_l}$  estimators are no longer independent and therefore  $V = \sum V_l/N_l$  type decomposition is no longer possible. The authors of [BJL<sup>+</sup>15] give an alternative argument to find (approximately) optimal  $N_l$ . They define

$$V_l := \left\| \frac{Z_{l-1}}{Z_l} G_{l-1} - 1 \right\|_\infty^2.$$

They show that  $N_l \propto \sqrt{V_l/C_l}$  is asymptotically optimal as the target error  $\varepsilon \rightarrow 0$  [BJL<sup>+</sup>15].

We now end this section by briefly discussing one of the important theoretical results from Beskos et.al. [BJL<sup>+</sup>15]. We denote  $h_l$  as the size of an element within the mesh, so that as the mesh resolution increases,  $h_l$  decreases. One can also think of this as a coarseness index of the mesh. To develop theory, one needs to characterize the asymptotic behaviour of bias, variance ( $V_l$ ) and complexity of forward map evaluation ( $C_l$ ):

1. Bias  $\sim \mathcal{O}(h_L^\alpha)$
2.  $C_l \sim \mathcal{O}(h_l^{-\zeta})$
3.  $V_l \sim \mathcal{O}(h_l^\beta)$

We assume that these quantities have the following relation :  $\alpha = \beta/2$  and  $\zeta \leq \beta$ . In practice, Bias and complexity results are usually available in the literature for many practically relevant models. On the other hand, the variance term depends on the specific Monte Carlo method and has to be established separately. We further assume the following;

1. There exists  $0 < \underline{C} < \overline{C} < +\infty$  s.t.

$$\sup_l \sup_u G_l(u) \leq \overline{C}$$

$$\inf_l \inf_u G_l(u) \geq \underline{C}.$$

2. There exists  $\rho \in (0, 1)$  s.t. for any  $l, A, (u, v) \in E^2$ , where  $E$  is the state space;

$$\int_A K_l(u, du') \geq \rho \int_A K_l(v, dv').$$

We can now state the theorem,

**Theorem 19.** *Assuming the aforementioned statements, in order to obtain an error of  $\mathcal{O}(\epsilon^2)$  (r.m.s. of  $\mathcal{O}(\epsilon)$ ), we only require a computational cost of  $\mathcal{O}(\epsilon^{-2})$ .*

Let us compare this result to the uni-resolution case, where we apply SMC to the finest resolution. Take  $\alpha = 2, \beta = 4, \zeta = 2$  With the aforementioned assumptions, we have:

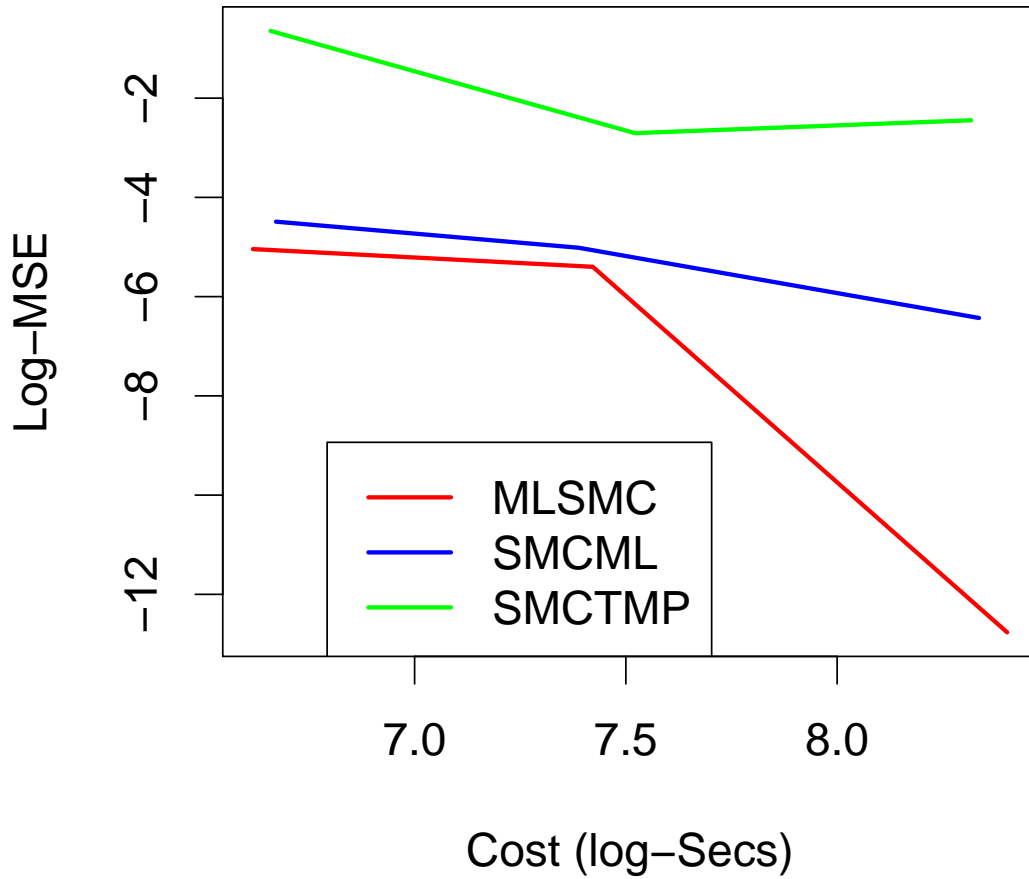
$$MSE \approx c_1 N^{-1} + c_2 h_L^{2\alpha}$$

which corresponds to the same variance, bias-squared decomposition as before. To get a r.m.s. error of  $\mathcal{O}(\epsilon)$ , we need  $N \sim \mathcal{O}(\epsilon^{-2})$  and  $h_L \sim \mathcal{O}(\epsilon^{0.5})$ . Therefore, the computational cost  $N \times C_L$  becomes  $\mathcal{O}(\epsilon^{-4})$ . It is interesting to note that, these complexity estimates agree with that of Hoang et.al. [HSS12] in the context of multi-resolution independence Samplers for the same elliptic problem. In our experience, the practicality of the independence samplers is limited due to the slow mixing of this proposal, however these results indicate that there is a family of multi-resolution Monte Carlo methods with the same asymptotic performance. Finally, we note that independence sampler is a special case of MCMC methods with dimensionality-independent mixing rates, as discussed in Section 2.1.3.

## 4.5 Results

Beskos et.al.[BJL<sup>+</sup>15] demonstrate the performance of their methodology and empirical validity of their theory using a simplified 1D problem. In this work, we implemented this method for the general 1D, 2D and 3D cases for both source and permeability inversion. The forcing inversion case, in particular, forms a basis to check the theoretically established rates. This is done in Figure 4.1. In this Figure, MLSMC denotes the method discussed in this chapter. On the other hand, SMCML is the SMC estimator where we use only the last

level of the particle system, in essence treating the level sequence as a bridging sequence. Finally, SMCTMP is the method described in the previous chapter. The theory indicates that SMCML and SMCTMP should have the same slope, while MLSMC having a smaller one. This result shows an agreement with the established theory. The smaller MSE per cost of SMCML compared to SMCTMP is due to the more computationally efficient construction of bridging densities.



**Figure 4.1:** log-Cost vs. log-Error plot for three methods

It is worth noting that our implementation is fully parallelized using MPI libraries as before. Additionally, we opted to use the adaptive SMC method developed in the last chapter, to

evolve the particle system from prior to the first level. In our observations, this approach performs much better than starting with a very coarse mesh resolution and evolving the system from prior in one step. We believe this behaviour can be explained as follows; the posterior corresponding to coarse resolutions are too different from those corresponding to fine resolutions (i.e. the bias component is too big), and this adversely effects the evolution of the particle system, in other words the low resolution posteriors are bad proposals for consecutive levels. This means, in applications, it is important to correctly select a good starting resolution; too coarse would be detrimental for inference, while too fine would eliminate any possible performance gain. The use of our adaptive SMC method alleviates this issue, and is one of our contributions.

Finally, the run time of this algorithm is around one-tenth to one-twentieth of uni-resolution SMC for the experiments considered. This gives hope for the practical usage of multi-resolution SMC in practical inverse problems. We should note, however, that the performance gap between Tikhonov-regularization based methods and multi-resolution SMC is still significant.

# CHAPTER 5

## Conclusion

Despite the pessimistic results of Bickel et.al. [BLB08, BBL08] about particle filters (and sequential Monte Carlo samplers in general) in high dimensional problems, recent theoretical [BCJ14, BCJW11, RvH15] and methodological advances [KBJ13] show the existence of a large class of problems that can be efficiently solved. To explore further the potential of sequential Monte Carlo samplers in important high dimensional inference problems, we focused on the Groundwater-Flow inverse problem. The high computational cost of likelihood computations and the existence of a natural sequence of increasingly accurate approximations make this problem an ideal test case. The recent research on dimensionality-robust Markov Chain Monte Carlo methods [BS09] also plays a role in our constructions.

Our first approach investigates the effects of novel and improved adaptations of various parameters associated with a SMC sampler using tempered posteriors as bridging densities. This approach allocates computational resources to match the acceptance rate of underlying MCMC steps to a given level. As a result, the highly tempered (“hot”) densities require less computational budget, so that increased effort can be put to more challenging target densities. As a side effect, this additional adaptivity makes the methodology a lot easier to use for the end-user, since very few parameters are needed to be manually specified, as opposed to the typical implementation.

The second approach utilizes the natural sequence of increasingly finer meshes in Finite Element approximations. The telescoping sum decomposition of the estimator forms the basis of this approach. By estimating each term in this sum separately, and using optimal number of particles, we can significantly reduce the estimation variance while keeping the

---

bias the same. Our adaptive SMC can also be employed here as an initial estimator. We observed that this multi-resolution approach yields significant efficiency gains and provides hope for the practical use of such an approach.

## References

- [AS09] Simon R Arridge and John C Schotland. Optical tomography: forward and inverse problems. *Inverse Problems*, 25(12):123010, 2009.
- [BBC85] J.V. Beck, B. Blackwell, and C.R.S. Clair. *Inverse Heat Conduction: Ill-Posed Problems*. Wiley-Interscience publication. Wiley, 1985.
- [BBG<sup>+</sup>11] L. Biegler, G. Biros, O. Ghattas, M. Heinkenschloss, D. Keyes, B. Mallick, L. Tenorio, B. Waanders, K. Willcox, and Y. Marzouk. *Large-Scale Inverse Problems and Quantification of Uncertainty*. Wiley Series in Computational Statistics. Wiley, 2011.
- [BBL08] Thomas Bengtsson, Peter Bickel, and Bo Li. *Curse-of-dimensionality revisited: Collapse of the particle filter in very large scale systems*, volume Volume 2 of *Collections*, pages 316–334. Institute of Mathematical Statistics, Beachwood, Ohio, USA, 2008.
- [BCJ14] Alexandros Beskos, Dan Crisan, and Ajay Jasra. On the stability of sequential monte carlo methods in high dimensions. *Ann. Appl. Probab.*, 24(4):1396–1445, 08 2014.
- [BCJW11] Alexandros Beskos, Dan Crisan, Ajay Jasra, and Nick Whiteley. Error bounds and normalizing constants for sequential monte carlo in high dimensions. *arXiv preprint arXiv:1112.1544*, 2011.
- [BGLR15] M. Banterle, C. Grazian, A. Lee, and C. P. Robert. Accelerating Metropolis-Hastings algorithms by Delayed Acceptance. *ArXiv e-prints*, March 2015.
- [BJKT] A Beskos, A Jasra, N Kantas, and A Thiery. On the convergence of adaptive sequential monte carlo methods.



- [BJL<sup>+</sup>15] Alexandros Beskos, Ajay Jasra, Kody Law, Raul Tempone, and Yan Zhou. Multilevel sequential monte carlo samplers. *arXiv preprint arXiv:1503.07259*, 2015.
- [BJMS15] Alexandros Beskos, Ajay Jasra, Ege A. Muzaffer, and Andrew M. Stuart. Sequential monte carlo methods for bayesian elliptic inverse problems. *Statistics and Computing*, 25(4):727–737, July 2015.
- [BLB08] Peter Bickel, Bo Li, and Thomas Bengtsson. *Sharp failure rates for the bootstrap particle filter in high dimensions*, volume Volume 3 of *Collections*, pages 318–329. Institute of Mathematical Statistics, Beachwood, Ohio, USA, 2008.
- [BPR07] Frank Bauer, Sergei Pereverzev, and Lorenzo Rosasco. On regularization algorithms in learning theory. *Journal of Complexity*, 23(1):52 – 72, 2007.
- [BRS09] Alexandros Beskos, Gareth Roberts, and Andrew Stuart. Optimal scalings for local metropolis–hastings chains on nonproduct targets in high dimensions. *Ann. Appl. Probab.*, 19(3):863–898, 06 2009.
- [BS09] Alexandros Beskos and Andrew Stuart. Mcmc methods for sampling function space. *Invited Lectures, Sixth International Congress on Industrial and Applied Mathematics, ICIAM07, European Mathematical Society*, pages 337–364, 2009.
- [BT] V. Bally and D. Talay. The law of the euler scheme for stochastic differential equations. *Probability Theory and Related Fields*, 104(1):43–60.
- [BTBG<sup>+</sup>12] T. Bui-Thanh, C. Burstedde, O. Ghattas, J. Martin, G. Stadler, and L. C. Wilcox. Extreme-scale uq for bayesian inverse problems governed by pdes. In *High Performance Computing, Networking, Storage and Analysis (SC), 2012 International Conference for*, pages 1–11, Nov 2012.
- [CF05] J. Andr s Christen and Colin Fox. Markov chain monte carlo using an approximation. *Journal of Computational and Graphical Statistics*, 14(4):795–810, 2005.

- [Cho04] Nicolas Chopin. Central limit theorem for sequential monte carlo methods and its application to bayesian inference. *Ann. Statist.*, 32(6):2385–2411, 12 2004.
- [CMW15] T. Cui, Y. M. Marzouk, and K. E. Willcox. Data-driven model reduction for the Bayesian solution of inverse problems. *International Journal for Numerical Methods in Engineering*, 102:966–990, May 2015.
- [CRSW13] S. L. Cotter, G. O. Roberts, A. M. Stuart, and D. White. Mcmc methods for functions: Modifying old algorithms to make them faster. *Statist. Sci.*, 28(3):424–446, 08 2013.
- [CRY09] Radu V. Craiu, Jeffrey Rosenthal, and Chao Yang. Learn from thy neighbor: Parallel-chain and regional adaptive mcmc. *Journal of the American Statistical Association*, 104(488):1454–1466, 2009.
- [DJ09] Arnaud Doucet and Adam M Johansen. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of Nonlinear Filtering*, 12(656-704):3, 2009.
- [DM04] P. Del Moral. *Feynman-Kac Formulae: Genealogical and Interacting Particle Systems with Applications*. Probability and Its Applications. Springer-Verlag New York, 2004.
- [DMDJ06] Pierre Del Moral, Arnaud Doucet, and Ajay Jasra. Sequential monte carlo samplers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3):411–436, 2006.
- [DS13] M. Dashti and A. M. Stuart. The Bayesian Approach To Inverse Problems. *ArXiv e-prints*, February 2013.
- [Dur] Ricardo G. Duran. Galerkin approximations and finite element methods. [http://mate.dm.uba.ar/~rduran/class\\_notes/fem.pdf](http://mate.dm.uba.ar/~rduran/class_notes/fem.pdf).
- [FN97] Colin Fox and Geoff Nicholls. Sampling conductivity images via mcmc. In *University of Leeds*, pages 91–100, 1997.

- [FWA<sup>+</sup>11] H. P. Flath, L. C. Wilcox, V. AkÅgelik, J. Hill, B. van Bloemen Waanders, and O. Ghattas. Fast algorithms for bayesian uncertainty quantification in large-scale linear inverse problems based on low-rank partial hessian approximations. *SIAM Journal on Scientific Computing*, 33(1):407–432, 2011.
- [gil08] Multilevel monte carlo path simulation. *Operations Research*, 56(3):607–617, 2008.
- [GT95] Charles J. Geyer and Elizabeth A. Thompson. Annealing markov chain monte carlo with applications to ancestral inference. *Journal of the American Statistical Association*, 90(431):909–920, 1995.
- [Gul05] Ramesh M. Gulrajani. *Modeling and Imaging of Bioelectrical Activity: Principles and Applications*, chapter The Forward Problem of Electrocardiography: Theoretical Underpinnings and Applications, pages 43–79. Springer US, Boston, MA, 2005.
- [HD12] T. P. Hill and M. Dall’Aglio. Bayesian Posteriors Without Bayes’ Theorem. *ArXiv e-prints*, March 2012.
- [HSS12] V. H. Hoang, C. Schwab, and A. M. Stuart. Complexity Analysis of Accelerated MCMC Methods for Bayesian Inversion. *ArXiv e-prints*, July 2012.
- [HST01] Heikki Haario, Eero Saksman, and Johanna Tamminen. An adaptive metropolis algorithm. *Bernoulli*, 7(2):223–242, 04 2001.
- [IJ14] K. Ito and B. Jin. *Inverse Problems: Tikhonov Theory and Algorithms*. Series on Applied Mathematics. World Scientific Publishing Company Pte Limited, 2014.
- [Ing97] Gabriele Inglese. An inverse problem in corrosion detection. *Inverse Problems*, 13(4):977, 1997.
- [JSDT11] AJAY JASRA, DAVID A. STEPHENS, ARNAUD DOUCET, and THEODOROS TSAGARIS. Inference for li£jvy-driven stochastic volatility

- models via adaptive sequential monte carlo. *Scandinavian Journal of Statistics*, 38(1):1–22, 2011.
- [JSH07] Ajay Jasra, David A. Stephens, and Christopher C. Holmes. On population-based simulation for static inference. *Statistics and Computing*, 17(3):263–279, 2007.
- [KBJ13] N. Kantas, A. Beskos, and A. Jasra. Sequential Monte Carlo Methods for High-Dimensional Inverse Problems: A case study for the Navier-Stokes equations. *ArXiv e-prints*, July 2013.
- [KPSC06] Benjamin S. Kirk, John W. Peterson, Roy H. Stogner, and Graham F. Carey. libmesh : a c++ library for parallel adaptive mesh refinement/coarsening simulations. *Engineering with Computers*, 22(3):237–254, 2006.
- [KS15] K. Konakli and B. Sudret. Polynomial meta-models with canonical low-rank approximations: numerical insights and comparison to sparse polynomial chaos expansions. *ArXiv e-prints*, November 2015.
- [LW15] A. Lee and N. Whiteley. Variance estimation and allocation in the particle filter. *ArXiv e-prints*, September 2015.
- [MN09] Youssef M. Marzouk and Habib N. Najm. Dimensionality reduction and polynomial chaos acceleration of bayesian inference in inverse problems. *Journal of Computational Physics*, 228(6):1862 – 1902, 2009.
- [MWBG12] James Martin, Lucas C. Wilcox, Carsten Burstedde, and Omar Ghattas. A stochastic newton mcmc method for large-scale statistical inverse problems with application to seismic inversion. *SIAM Journal on Scientific Computing*, 34(3):A1460–A1487, 2012.
- [Nat01] F. Natterer. *The Mathematics of Computerized Tomography*. Society for Industrial and Applied Mathematics, 2001.
- [Neu98] Arnold Neumaier. Solving ill-conditioned and singular linear systems: A tutorial on regularization. *SIAM Review*, 40(3):636–666, 1998.

- [Pes73] P. H. Peskun. Optimum monte carlo sampling using markov chains. *Biometrika*, 60:607–612, 1973.
- [RC05] Christian P. Robert and George Casella. *Monte Carlo Statistical Methods (Springer Texts in Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2005.
- [RfM03] E. Roeckner and Max-Planck-Institut für Meteorologie. *The Atmospheric General Circulation Model ECHAM5: Part 1 : Model Description*. Max-Planck-Institut für Meteorologie, 2003.
- [Ron08] Luca Rondi. On the regularization of the inverse conductivity problem with discontinuous conductivities. *Inverse Problems and Imaging*, 2(3):397–409, 2008.
- [RR04] Gareth O. Roberts and Jeffrey S. Rosenthal. General state space markov chains and mcmc algorithms. *Probab. Surveys*, 1:20–71, 2004.
- [RvH15] Patrick Rebeschini and Ramon van Handel. Can local particle filters beat the curse of dimensionality? *Ann. Appl. Probab.*, 25(5):2809–2866, 10 2015.
- [SOL<sup>+</sup>12] Antti Solonen, Pirkka Ollinaho, Marko Laine, Heikki Haario, Johanna Tamminen, and Heikki Järvinen. Efficient mcmc for climate model parameter estimation: Parallel adaptive chains and early rejection. *Bayesian Anal.*, 7(3):715–736, 09 2012.
- [SS02] M. Schatzman and M. Schatzman. *Numerical Analysis: A Mathematical Introduction*. Clarendon Press, 2002.
- [SV82] L. A. Shepp and Y. Vardi. Maximum likelihood reconstruction for emission tomography. *IEEE Transactions on Medical Imaging*, 1(2):113–122, Oct 1982.
- [Ten01] Luis Tenorio. Statistical regularization of inverse problems. *SIAM Review*, 43(2):347–366, 2001.
- [Vol13] Sebastian J Vollmer. Posterior consistency for bayesian inverse problems through stability and regression results. *Inverse Problems*, 29(12):125011, 2013.

- 
- [WA11] Kun Wang and Mark A. Anastasio. *Handbook of Mathematical Methods in Imaging*, chapter Photoacoustic and Thermoacoustic Tomography: Image Formation Principles, pages 781–815. Springer New York, New York, NY, 2011.
- [Zhd93] Michael S. Zhdanov. Tutorial-Regularization in inversion theory.pdf, 1993.